



ANALYSIS OF GLOBAL AND LOCAL OPTIMA OF REGULARIZED QUANTILE REGRESSION IN HIGH DIMENSIONS: A SUBGRADIENT APPROACH

LAN WANG 
University of Miami

XUMING HE 
University of Michigan

Regularized quantile regression (QR) is a useful technique for analyzing heterogeneous data under potentially heavy-tailed error contamination in high dimensions. This paper provides a new analysis of the estimation/prediction error bounds of the global solution of L_1 -regularized QR (QR-LASSO) and the local solutions of nonconvex regularized QR (QR-NCP) when the number of covariates is greater than the sample size. Our results build upon and significantly generalize the earlier work in the literature. For certain heavy-tailed error distributions and a general class of design matrices, the least-squares-based LASSO cannot achieve the near-oracle rate derived under the normality assumption no matter the choice of the tuning parameter. In contrast, we establish that QR-LASSO achieves the near-oracle estimation error rate for a broad class of models under conditions weaker than those in the literature. For QR-NCP, we establish the novel results that all local optima within a feasible region have desirable estimation accuracy. Our analysis applies to not just the hard sparsity setting commonly used in the literature, but also to the soft sparsity setting which permits many small coefficients. Our approach relies on a unified characterization of the global/local solutions of regularized QR via subgradients using a generalized Karush–Kuhn–Tucker condition. The theory of the paper establishes a key property of the subdifferential of the quantile loss function in high dimensions, which is of independent interest for analyzing other high-dimensional nonsmooth problems.

1. INTRODUCTION

The semiparametric technique of quantile regression (QR) provides a useful alternative to least-squares regression and has been widely applied to data analysis

Wang and He's research is partly supported by NSF FRGMS-1952373. The authors are grateful to the Co-Editor and two anonymous referees, whose comments have helped to significantly improve the paper. They also thank Dr. Alexander Giessing for his helpful comments and Dr. Yunan Wu for her latex help on an earlier draft of the paper. Part of the results developed in this paper were made available as an earlier technical report (Wang, 2019). Address correspondence to Lan Wang, Department of Management Science, University of Miami, Coral Gables, FL 33146, USA; e-mail: lanwang@mbs.miami.edu.

in economics and finance since its introduction in the seminal paper of Koenker and Bassett (1978). For example, a low quantile of the return distribution of an investment portfolio provides an assessment of risk commonly known as Value at Risk. Buchinsky (1994), Chamberlain (1994), Buchinsky (1998), Abadie, Angrist, and Imbens (2002), Horowitz and Spokoiny (2002), Angrist, Chernozhukov, and Fernández-Val (2006), Firpo, Fortin, and Lemieux (2009), Galvao, Lamarche, and Lima (2013), Arellano and Bonhomme (2017), and Graham et al. (2018), among others, employed QR to study the wage distribution. See also Linton and Whang (2004), Horowitz and Lee (2005), Koenker and Xiao (2006), Chernozhukov and Fernández-Val (2011), Chernozhukov et al. (2013), Fitzenberger, Koenker, and Machado (2013), Su and Hoshino (2016), Koenker (2017), and Koenker et al. (2017) for other interesting applications of QR in economics. QR helps characterize the entire conditional distribution and often leads to insightful discoveries that would otherwise be imperceptible. It also has the appealing property of being robust to heavy-tailed error distributions. By contrast, L_1 -regularized least-squares regression, to be called LS-LASSO throughout the paper (Tibshirani, 1996), is known to be vulnerable to heavy-tailed errors.

Let Y be a random variable, and let $X = (x_1, \dots, x_p)^T$ be a p -dimensional vector of covariates. A linear QR model takes the form

$$Y = X^T \beta^* + \epsilon, \quad P(\epsilon \leq 0|X) = \tau, \quad \text{for some } 0 < \tau < 1, \quad (1)$$

where the error distribution of ϵ is generally heteroskedastic, and $\beta^* = (\beta_1^*, \beta_2^*, \dots, \beta_p^*)^T$ is the unknown parameter vector. In this formulation, both ϵ and β may depend on the quantile level τ of interest, but we ignore such dependence in notation for simplicity. Model (1) implies that $Q_{Y|X}(\tau) = X^T \beta^*$, where $Q_{Y|X}(\tau) = \inf\{t : F_{Y|X}(t) \geq \tau\}$ is the τ th conditional quantile of Y given X . We are interested in estimating β^* in the setting where the number of covariates p is much larger than the sample size n .

This paper develops a useful technique to study high-dimensional QR in a general framework under a set of lean assumptions. Our theory relies on establishing a key *restricted strong convexity* (RSC) property of the subdifferential of the quantile loss function in high dimensions. Let $S_n(\beta)$ be any subgradient of the sample quantile loss function (see more details in Section 2.3 of the main paper) and denote $\Delta = \beta - \beta^*$. Theorem 1 of Section 2.2 shows that there exist some positive constants α^* and c^* that do not depend on n or p such that

$$\langle S_n(\beta) - S_n(\beta^*), \Delta \rangle \geq \alpha^* \|\Delta\|_2^2 - c^* \sqrt{\frac{\log p}{n}} \|\Delta\|_1,$$

uniformly on $\{\|\Delta\|_2 \leq 1\} \cap \mathbb{C}$, with high probability, where the set \mathbb{C} , to be made clearer later in the paper, depends on the specific regularization method and the sparsity pattern of the true parameter β^* .

The subgradient approach leads to a unified analysis of both the global solution of the L_1 -regularized QR (QR-LASSO) and the local solutions of nonconvex regularized QR (QR-NCP) under a set of mild assumptions while allowing for

a more general sparsity pattern of β^* . We include both QR-LASSO and QR-NCP in the analysis because both types of penalty functions are of substantial interest in the literature and in practice. QR-LASSO is computationally convenient due to the convexity of the objective function, whereas QR-NCP helps alleviate the estimation bias due to the potential over-penalization of the L_1 -penalty. For regularized least-squares regression and generalized linear models in high dimensions, an equivalent property of RSC has been shown to play a fundamental role in theoretical analysis. However, such results are available only for differentiable loss functions. The gradient function of the quantile loss is not even Lipschitz continuous, which leads to substantial technical challenges. Our proof involves a novel construction of a Lipschitz continuous lower bound and the use of modern empirical processes techniques.

In the classical setting where the number of covariates is not large, Wang, Li, and Jiang (2007), Li and Zhu (2008), Zou and Yuan (2008), Wu and Liu (2009), Shows, Lu, and Zhang (2010), Kai, Li, and Zou (2011), Wagener, Volgushev, and Dette (2012), Wang, Zhou, and Li (2013a), and Chen et al. (2019a), among others, investigated regularized QR for variable selection. Several authors have recently investigated QR in high dimensions. Belloni and Chernozhukov (2011) was among the first to rigorously establish estimation error bounds for QR-LASSO (see also Kato, 2011; Wang, 2013). More recently, Park, He, and Zhou (2017) investigated multiple QR with high-dimensional covariates, Lee et al. (2018) studied high-dimensional QR-LASSO with a change point, Harding and Lamarche (2018) investigated QR-LASSO for panel data, and Chen, Liu, and Zhang (2019b) explored QR for big data under memory constraint. Moreover, adaptively weighted QR-LASSO or QR-NCP has been considered for better variable selection performance in various settings (see Bradic, Fan, and Wang, 2011; Wang, Wu, and Li, 2012; Fan, Fan, and Barut, 2014; Zheng, Peng, and He, 2015, among others). High-dimensional semiparametric QR has been investigated by Tang et al. (2013), Sherwood and Wang (2016), Zhong et al. (2016), Fan and Lian (2018), Lv et al. (2018), and Honda, Ing, and Wu (2019), among others.

Inspired by the recent work, this paper makes several new contributions to the fundamental theory of QR in the regime where the number of covariates p can grow at an exponential rate of the sample size n .

- We show that QR-LASSO enjoys the near-oracle estimation error rate under a set of lean assumptions. Our theory permits a rich class of error distributions as well as a general class of random design matrices without requiring the *nonlinear eigenvalue condition*. The estimation error rates are established under not only the popular hard sparsity setting, but also a more relaxed soft sparsity assumption which permits many covariates to have small effects.
- For QR-NCP, we show that all local minima within a feasible region have desirable error bounds and achieve the minimax error rate of estimation. These new results fill an important theoretical gap in the literature, because the global minimizers for nonconvex objective functions are not always numerically obtained or verifiable in practice.

- We derive the quantile prediction error rate by a general characterization of the prediction error based on subgradients.

Our results demonstrate that QR-LASSO enjoys near-oracle estimation error rates for a much richer class of error distributions than LS-LASSO does. The theory relies on conditions generally weaker than those in the current literature for LS-LASSO and QR-LASSO. Our analysis of the local minima for QR-NCP is new. Computation of the global solution of QR-NCP is infeasible in high dimensions. On the other hand, the empirical results in the literature demonstrate that the local solutions (obtained by different algorithms) of QR-NCP often significantly reduce the bias of QR-LASSO. The existing theory of QR-NCP has been focused on the existence of a local solution with good statistical properties. One main contribution of our paper is to fill in the gap of the theory by establishing that any solution satisfying the first-order condition (including global minimum) within a given radius of the true value has the desirable statistical accuracy. This result has important implications for the use of QR-NCP as an estimator for the QR coefficient or as an initial value for inference (see Section 6 for more discussions). The results substantially generalize those in Loh and Wainwright (2015), Mei, Bai, and Montanari (2018), and Elsener and van de Geer (2018) on the properties of local minima for differentiable loss functions. It is worth noting that the statistical properties of local solutions are of broader interest. Even in low dimensions, algorithms for many nonlinear problems (e.g., nonlinear generalized method of moments (GMM)) only guarantee first-order solutions (e.g., stationary point). Moreover, our subgradient approach is different from the techniques commonly used in the literature of high-dimensional QR and is of significant independent interest. The proposed technique is applicable to a large class of high-dimensional nonsmooth problems, for instance, classification based on the hinge loss.

The rest of the paper is organized as follows. Section 2 introduces the background and a general characterization of the solution of regularized QR based on subgradients. Section 3 presents the main theory on the estimation error bounds for QR-LASSO and QR-SCAD in ultra-high dimensions. Section 4 studies the quantile prediction error bounds. Section 5 reports results from a Monte Carlo study. Section 6 concludes the paper with additional remarks. Appendixes A–D contain detailed technical arguments.

1.1. Notation

For any vector $v = (v_1, \dots, v_p)^T \in \mathbb{R}^p$, $\|v\|_2 = \sqrt{\sum_{i=1}^p v_i^2}$ denotes its L_2 -norm, $\|v\|_1 = \sum_{i=1}^p |v_i|$ denotes its L_1 -norm, and $\|v\|_\infty = \max_{1 \leq i \leq p} |v_i|$ denotes its L_∞ -norm. Given an arbitrary index set $S \subseteq \{1, \dots, p\}$, v_S denotes the subvector of v containing the elements whose indices are in S , v_{S^c} denotes the subvector of v containing the elements whose indices are not in S , and $|S|$ denotes the cardinality of the index set S . For a real symmetric matrix M , $\lambda_{\max}(M)$ denotes its largest eigenvalue. For sequences of real positive numbers a_n and b_n , $a_n \sim b_n$ means $c_1 \leq$

$a_n/b_n \leq c_2$ for some positive constants c_1 and c_2 . The covariates $X_i = (x_{i1}, \dots, x_{ip})^T$ are independent p -dimensional sub-Gaussian random vectors with variance proxy ζ_0^2 , that is, $\forall v \in \mathbb{R}^p, \forall t \in \mathbb{R}, \mathbb{E}\{\exp(tX_i^T v)\} \leq \exp\{\zeta_0^2 t^2 \|v\|_2^2/2\}$, where ζ_0 is a positive constant. In this paper, a number is referred to as “a constant” if it does not depend on (n, p) , but it is allowed to depend on the underlying probability distributions of X_i and ϵ_i . The indicator function of an event A , denoted by $I(A)$, takes the value one if A occurs and zero otherwise. We often write $\frac{1}{p}$ as $\exp(-\log p)$ to emphasize that this term converges to zero at an exponential rate of $\log p$. The sign function $\text{sgn}(t) = 1$ if $t > 0$; -1 if $t < 0$; and takes its value in $[-1, 1]$ if $t = 0$.

2. PRELIMINARIES

In this section, we briefly review regularized QR with convex and nonconvex penalties, under the hard sparsity or soft sparsity assumption. We then elaborate on how to use subgradients to characterize the regularized QR solutions, global or local, with a general penalty function.

2.1. Background

Consider a random sample $\{Y_i, X_i\}_{i=1}^n$ satisfying model (1), where $X_i = (x_{i1}, \dots, x_{ip})'$. To explicitly incorporate the intercept term, we assume $x_{i1} = 1$ for all i , and correspondingly $\beta^* = (\beta_1^*, \beta_-^{*T})^T$ where $\beta_-^* = (\beta_2^*, \dots, \beta_p^*)^T$.

To avoid overfitting in the setting $p \gg n$, we consider estimating β^* by regularized QR:

$$\hat{\beta} \equiv \hat{\beta}(\lambda) \equiv \arg \min_{\beta \in \mathbb{R}^p} \left\{ n^{-1} \sum_{i=1}^n \rho_\tau(Y_i - X_i^T \beta) + \sum_{j=2}^p p_{\lambda_n}(|\beta_j|) \right\}, \quad (2)$$

where $\lambda = \lambda_n$ is positive tuning parameter controlling the complexity of the solution, $\rho_\tau(u) = u\{\tau - I(u < 0)\}$ is the quantile loss function, $\beta = (\beta_1, \beta_-^T)^T \in \mathbb{R}^p$, and $p_{\lambda_n}(\cdot)$ denotes a penalty function with tuning parameter λ_n . A larger value of λ_n encourages sparser solutions.

QR-LASSO adopts the popular L_1 -penalty for which $p_{\lambda_n}(|\beta_j|) = \lambda_n |\beta_j|$. It is computationally convenient due to its convex structure. Alternatively, one may use a nonconvex penalty function, which can alleviate the bias associated with the L_1 -penalty. Two popular choices of nonconvex penalty functions are the SCAD penalty function (Fan and Li, 2001) and the MCP penalty function (Zhang, 2010). The SCAD penalty function is defined as

$$\begin{aligned} p_\lambda(|\beta|) &= \lambda |\beta| I(|\beta| < \lambda) + \frac{a\lambda |\beta| - (\beta^2 + \lambda^2)/2}{a-1} I(\lambda \leq |\beta| \leq a\lambda) \\ &\quad + \frac{(a+1)\lambda^2}{2} I(|\beta| > a\lambda), \text{ for some } a > 2; \end{aligned}$$

whereas the MCP function has the form

$$p_\lambda(|\beta|) = \lambda \left(|\beta| - \frac{\beta^2}{2a\lambda} \right) I(|\beta| < a\lambda) + \frac{a\lambda^2}{2} I(|\beta| \geq a\lambda), \quad \text{for some } a > 1.$$

Most of the existing literature on high-dimensional QR assumes that β^* satisfies a hard sparsity constraint, especially

$$\beta^* \in \mathbb{B}_0(s) = \left\{ \beta \in \mathbb{R}^p : \sum_{j=2}^p I(\beta_j \neq 0) \leq s-1 \right\}, \quad (3)$$

for some positive constant $1 \leq s \ll n$, where $I(\cdot)$ denotes the indicator function. Hence, $\|\beta^*\|_0 \leq s$ and most of the components in β^* are exactly zero. As is customary for a regression model in the classical setting, the intercept term is always included in the model. The sparsity constraints are thus imposed on the slope components of β^* . Note that this is a subtle difference from high-dimensional least-squares regression where the intercept is generally taken as zero, which can be done for mean regression without loss of generality by centering both the response variable and the regressors. However, such a simplification does not carry over to QR.

The hard sparsity constraint may be overly restrictive for some applications which involve many weak signals, rather than just a few strong signals. This paper also considers a more relaxed sparsity constraint, which allows β^* to have many smallish nonzero coefficients. More specifically, the soft sparsity constraint assumes

$$\beta^* \in \mathbb{B}_1(R) = \left\{ \beta \in \mathbb{R}^p : \sum_{j=2}^p |\beta_j| \leq R \right\} \quad (4)$$

for some positive number R , which may depend on the sample size. In (4), instead of L_1 -norm, we may also use the L_q -norm for some $0 < q < 1$. The results of the paper would still hold under minor modifications. It is worth noting that both $\mathbb{B}_0(s)$ and $\mathbb{B}_1(R)$ may depend on the quantile level τ of interest. Discussions on the identification of the population parameter β_0 in the high-dimensional setting are given in Appendix A.

2.2. Characterizing the Solutions for Regularized Quantile Regression

We now present a unified characterization of the regularized QR estimators, including the global solution of QR-LASSO and the local solutions of QR-NCP, based on a generalized Karush–Kuhn–Tucker (KKT) condition for the convex difference problem, characterized by subgradients.

A subgradient of a convex function $g(\beta)$ at β_1 is any vector $\xi \in \mathbb{R}^p$ such that $g(\beta_2) \geq g(\beta_1) + \xi^T(\beta_2 - \beta_1)$ for all β_2 . The subdifferential of $g(\beta)$ at β_1 , denoted by $\partial g(\beta_1)$, consists of all the subgradients of $g(\beta)$ at β_1 .

Let

$$Q_n(\beta) = n^{-1} \sum_{i=1}^n \rho_\tau(Y_i - X_i^T \beta) \quad (5)$$

denote the sample quantile loss function. One can show that the subdifferential $\partial Q_n(\beta)$ comprises vectors $S_n(\beta) = (S_{n1}(\beta), \dots, S_{np}(\beta))^T$, where, for $j = 1, \dots, p$,

$$\begin{aligned} S_{nj}(\beta) = & -\tau n^{-1} \sum_{i=1}^n x_{ij} \mathbf{I}(Y_i - X_i^T \beta > 0) \\ & + (1 - \tau) n^{-1} \sum_{i=1}^n x_{ij} \mathbf{I}(Y_i - X_i^T \beta < 0) - n^{-1} \sum_{i=1}^n x_{ij} v_i \end{aligned} \quad (6)$$

and

$$v_i \in \begin{cases} \{0\}, & \text{if } Y_i - X_i^T \beta \neq 0, \\ [\tau - 1, \tau], & \text{otherwise.} \end{cases}$$

Let $L_n(\beta) = Q_n(\beta) + \sum_{j=2}^p p_{\lambda_n}(|\beta_j|)$ be the regularized quantile objective function in (2). We observe that for a general class of penalty functions, $L_n(\beta)$ can be written as the difference of two convex functions in β :

$$L_n(\beta) = \tilde{L}_n(\beta) - H(\beta),$$

where $\tilde{L}_n(\beta) = Q_n(\beta) + \lambda \sum_{j=2}^p |\beta_j|$ and $H(\beta) = \sum_{j=2}^p h_\lambda(\beta_j)$, where $h_\lambda(\cdot)$ is differentiable. In the case of LASSO penalty, $h_\lambda(\beta_j) = 0$, for $j = 2, \dots, p$, and thus $\tilde{L}_n(\beta)$ coincides with $L_n(\beta)$. For the nonconvex SCAD penalty,

$$\begin{aligned} h_\lambda(\beta_j) = & \left[(\beta_j^2 - 2\lambda|\beta_j| + \lambda^2) / (2(a-1)) \right] \mathbf{I}(\lambda \leq |\beta_j| \leq a\lambda) \\ & + \left[\lambda|\beta_j| - (a+1)\lambda^2/2 \right] \mathbf{I}(|\beta_j| > a\lambda), \end{aligned}$$

whereas for the nonconvex MCP function,

$$h_\lambda(\beta_j) = \left[\beta_j^2 / (2a) \right] \mathbf{I}(|\beta_j| < a\lambda) + \left[\lambda|\beta_j| - a\lambda^2/2 \right] \mathbf{I}(|\beta_j| \geq a\lambda).$$

For the above convex difference optimization problem, an extension of the KKT condition was given in Tao and An, 1997, which implies that the solution $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^T$ of (2), global or local, satisfies the following necessary condition:

$$\nabla \tilde{L}_n(\hat{\beta}) - H'(\hat{\beta}) = 0, \quad (7)$$

where $\nabla \tilde{L}_n(\hat{\beta})$ denotes some (not necessarily any) subgradient in the subdifferential of \tilde{L}_n being evaluated at $\hat{\beta}$, and $H'(\hat{\beta}) = (0, h'_\lambda(\hat{\beta}_2), \dots, h'_\lambda(\hat{\beta}_p))^T$. Note that $\tilde{L}_n(\beta)$ is the sum of two convex functions. In this paper, we consider stationary points satisfying (7), which include the solutions given by popular algorithms.

LEMMA 1. Let $\widehat{\beta}$ be any stationary point of QR-LASSO or QR-NCP satisfying (7). Then there exists a subgradient $S_n(\beta) \in \partial Q_n(\beta)$ such that

$$S_n(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}) - H'(\widehat{\beta}) = 0, \quad (8)$$

where $\text{sgn}(\widehat{\beta}) = (0, \text{sgn}(\widehat{\beta}_2), \dots, \text{sgn}(\widehat{\beta}_p))^T$.

Consider next a particular subgradient \widetilde{S}_n in $\partial Q_n(\beta)$, given by

$$\widetilde{S}_n = n^{-1} \sum_{i=1}^n X_i \xi_i, \quad (9)$$

where $\xi_i = I(\epsilon_i < 0) - \tau$. For regularized QR, λ is usually selected such that the event

$$\Lambda_n = \{\lambda \geq c_0 \|\widetilde{S}_n\|_\infty\} \quad (10)$$

happens with high probability for some positive constant $c_0 > 1$. Such a penalty was also considered in Belloni and Chernozhukov (2011) for QR-LASSO (see also Kato (2011) for an extension to the group LASSO setting). The above choice of λ is motivated by the general principle of tuning parameter selection in regularized least-squares regression (Bickel, Ritov, and Tsybakov, 2009) and the KKT condition for a general convex difference problem (Tao and An, 1997). Following the choice for LS-LASSO (Bickel et al., 2009), we take $c_0 = 2$ in the subsequent analysis.

3. MAIN THEORY

In this section, we provide details on the theoretical properties of the regularized QR estimator $\widehat{\beta}$. For QR-LASSO, $\widehat{\beta}$ denotes the global solution defined in (2), which also satisfies (7). For QR-NCP, $\widehat{\beta}$ denotes any local solution satisfying (7).

3.1. Geometric Structure of Regularized Quantile Regression Estimators

Under the hard sparsity condition, the QR-LASSO estimator is known to lie in a cone-shaped set with high probability. Let $\widehat{v} = \widehat{\beta} - \beta^*$. Here, we go beyond the setting of the L_1 -penalty and hard sparsity to characterize the geometric structure of \widehat{v} .

Let $S_- = \{j : \beta_j^* \neq 0, 2 \leq j \leq p\}$ and $S = S_- \cup \{1\}$. Given a threshold $a > 0$, let $S_{-a} = \{j : |\beta_j^*| > a, 2 \leq j \leq p\}$. Let $S_a = S_{-a} \cup \{1\}$. The cardinality $\|S\|_0 = s$ denotes the sparsity size under the hard sparsity condition. Under the soft sparsity assumption, s can be much larger than n .

Let

$$\Gamma_H = \{v \in \mathbb{R}^p : \|v_{S^c}\|_1 \leq 3\|v_S\|_1\}, \quad (11)$$

$$\Gamma_W = \{v \in \mathbb{R}^p : \|v_{S_a^c}\|_1 \leq 3\|v_{S_a}\|_1 + 4\|\beta_{S_a^c}^*\|_1\}, \quad (12)$$

$$\tilde{\Gamma}_H = \{v \in \mathbb{R}^p : \|v_{A^c}\|_1 \leq 3\|v_A\|_1\}, \quad (13)$$

$$\tilde{\Gamma}_W = \{v \in \mathbb{R}^p : \|v_{S_a^c}\|_1 \leq 3\|v_{S_a}\|_1 + 2\|\beta_{S_a^c}^*\|_1\}, \quad (14)$$

where A is the index set corresponding to the s -largest (in absolute value) elements of v . It is observed that Γ_H and $\tilde{\Gamma}_H$ are cone-shaped, but Γ_W and $\tilde{\Gamma}_W$ are star-shaped. For example, if $v \in \Gamma_W$, then the whole line segment $\{tv | t \in (0, 1)\}$ is contained in Γ_W . The sets Γ_W and $\tilde{\Gamma}_W$ depend on a , but we omit the dependence in notation for simplicity.

We use $f_i(t)$ to denote the conditional probability density function of ϵ_i given X_i , $i = 1, \dots, n$. We also assume, without loss of generality, that the covariate $X_{-i} = (x_{i2}, \dots, x_{ip})^T$ is a $(p-1)$ -dimensional mean-zero random vector, and $\Sigma = E(X_i X_i^T)$ exists. Conditions (C1)–(C3) below constitute a set of basic assumptions for establishing the statistical properties of the regularized QR estimator in high dimensions.

Condition (C1). The conditional distribution of ϵ_i satisfies $P(\epsilon_i \leq 0 | X_i) = \tau$, $i = 1, \dots, n$. There exist positive constants m_0 and b_0 such that $\inf_{1 \leq i \leq n} f_i(t) \geq m_0 > 0$, for all $|t| \leq b_0$.

Condition (C2). The matrix Σ satisfies $\lambda_{\max}(\Sigma) \leq k_u < \infty$, and

$$v^T \Sigma v \geq m_1 \|v\|_2^2, \quad \text{for any } v \in \mathbb{C}, \quad (15)$$

for some constant m_1 , where for QR-LASSO, $\mathbb{C} = \Gamma_H$ under the hard sparsity assumption, and $\mathbb{C} = \Gamma_W$ under the soft sparsity assumption whereas for QR-NCP, $\mathbb{C} = \mathbb{R}^p$.

Condition (C3). Let $\hat{\sigma}_j^2 = n^{-1} \sum_{i=1}^n x_{ij}^2$. There exist a constant $m_x > 0$ and a positive sequence of numbers δ_n such that $P(\max_{1 \leq j \leq p} \hat{\sigma}_j^2 \leq m_x) \geq 1 - \delta_n$, where $\delta_n \rightarrow 0$ as $n \rightarrow \infty$.

Remark 1. Condition (C1) imposes regularity conditions on the random error distributions, which allow for heteroskedastic error distributions and do not require the existence any moment. The constant b_0 in (C1) may depend on the probability distribution of X_i , as described in Lemma C.2 of Appendix C. A large class of heavy-tailed error distributions, such as the Cauchy distribution, satisfy condition (C1). The restricted eigenvalue condition in (C2) is similar to those imposed for regularized least-squares regression. For QR-LASSO, (15) is exactly the same as the restricted eigenvalue condition for LS-LASSO as the restriction sets Γ_H and Γ_W have the same forms as those for LS-LASSO. For QR-NCP, the requirement of $\mathbb{C} = \mathbb{R}^p$ amounts to assuming $\lambda_{\min}(\Sigma) \geq m_1 > 0$. However, if we restrict our attention to a sparse local solution, then this can

be replaced by weaker sparse eigenvalue condition in Zhang (2010). Finally, condition (C3) is satisfied if the covariates have sub-Gaussian distributions. It can also be satisfied when some of the covariates do not have sub-Gaussian distributions. For example, if a small subset (say fixed size) of covariates only have finite second moments, whereas the others follow the sub-Gaussian distributions with bounded variance proxy, then (C3) still holds. Overall, the above set of conditions is similar to or weaker than that in the literature for high-dimensional QR. Some detailed comparisons are given in Remark 3 of Section 3.3.

LEMMA 2 (QR-LASSO). *Assume $\lambda = k_0\sqrt{\log p/n}$, where $k_0 \geq 4\sqrt{m_x}$ is a constant. Suppose conditions (C1) and (C3) are satisfied. Then, with probability at least $1 - \delta_n - 2\exp(-\log p)$, (i) $\hat{v} \in \Gamma_H$ under the hard sparsity assumption and (ii) $\hat{v} \in \Gamma_W$ under the soft sparsity assumption.*

For QR-LASSO, the geometric structure is a result of the convexity of the regularized quantile loss function. The first part of Lemma 2 under hard sparsity was observed in Belloni and Chernozhukov (2011), whereas the result under soft sparsity is new and is a generalization of Negahban et al. (2012). For QR-NCP, the geometric structure is less transparent. Instead, the structure is implicit in the derivation of the estimation error bound.

For QR-NCP, we have $\hat{v} \in \tilde{\Gamma}_H$ under the hard sparsity assumption and $\hat{v} \in \tilde{\Gamma}_W$ under the soft sparsity assumption with high probability. Due to the reliance on the conditions of later theorems, we refer to Corollary 1 in Section 3.3 for a full description of the results.

3.2. Properties of Subgradients in High Dimensions

We first state a useful property for the subgradient \tilde{S}_n defined in (9). The following lemma gives a high probability bound for its supremum norm.

LEMMA 3. *Suppose conditions (C1) and (C3) are satisfied. We have*

$$P(|\tilde{S}_n|_\infty \leq 2\sqrt{m_x \log p/n}) \geq 1 - \delta_n - 2\exp(-\log p).$$

The lemma suggests that the event Λ_n defined in (10) occurs with a high probability for an appropriate choice of λ at the rate $\sqrt{\log p/n}$.

Theorem 1 provides a core result for establishing error bounds for QR-LASSO and QR-NCP by showing that a type of restricted convexity condition holds with high probability for any subgradient in the subdifferential of the sample quantile loss function.

THEOREM 1 (Restricted strong convexity). *Suppose conditions (C1)–(C3) are satisfied. There exist some positive constants a^* , c^* , a_1 , and a_2 , such that for any*

subgradient $S_n \in \partial Q_n(\beta)$,

$$\langle S_n(\beta^* + \Delta) - S_n(\beta^*), \Delta \rangle \geq a^* \|\Delta\|_2^2 - c^* \sqrt{\frac{\log p}{n}} \|\Delta\|_1, \quad (16)$$

uniformly on $\{\|\Delta\|_2 \leq 1\} \cap \mathbb{C}$, holds with probability at least $1 - \delta_n - a_1 \exp(-a_2 \log p)$, where \mathbb{C} is defined in condition (C2) and $\delta_n \rightarrow 0$ is given in condition (C3).

Remark 2. This theorem guarantees an RSC condition on $\{\|\Delta\|_2 \leq 1\} \cap \mathbb{C}$. Lemma C.4 in Appendix C shows that a slightly weaker result holds uniformly on $\{\|\Delta\|_2 > 1\} \cap \mathbb{C}$. Specifically, with probability at least $1 - \delta_n - a_1 \exp(-a_2 \log p)$, $\langle S_n(\beta^* + \Delta) - S_n(\beta^*), \Delta \rangle \geq a^* \|\Delta\|_2 - c^* \sqrt{\frac{\log p}{n}} \|\Delta\|_1$, uniformly on $\{\|\Delta\|_2 > 1\} \cap \mathbb{C}$. Similar RSC-type conditions have been considered in the recent literature on high-dimensional M -regression (Negahban et al., 2012; Loh and Wainwright, 2015, among others). However, the existing literature has only considered smooth (second-order differential) loss functions. To the best of our knowledge, this is the first time the RSC condition is established for a nonsmooth loss function, which is more technically challenging due to the fact the gradient function is not even Lipschitz continuous. Our proof is based on a novel construction of a Lipschitz continuous lower bound and the application of advanced empirical process theory techniques (e.g., peeling). Our approach can be applied to other nonsmooth high-dimensional problems and is of interest beyond QR.

3.3. Estimation Error Bounds

This subsection derives the L_2 -error and L_1 -error bounds for the estimator $\hat{\beta}$ under both the hard sparsity assumption and the soft sparsity assumption. It is worth emphasizing that the results here are nonasymptotic in the sense that the error bounds hold for any (n, p) satisfying the stated conditions. The theory allows the number of covariates p to grow at an exponential rate of the sample size n , often called the ultra-high-dimensional setting. In the sequel, a^* , c^* , a_1 , and a_2 denote the positive constants given in Theorem 1.

THEOREM 2 (QR-LASSO). *Suppose conditions (C1)–(C3) are satisfied. Let $\lambda = k_0 \sqrt{\log p/n}$, where $k_0 \geq 4\sqrt{m_x}$ is a constant.*

- (i) *(Hard sparsity case) Assume β^* satisfies the hard sparsity assumption (3), with $n > (a_1^*)^2 s \log p$, where $a_1^* = 4(2k_0 + c^*)/a^*$. Then, with probability at least $1 - 4\delta_n - 4\exp(-\log p) - 2a_1 \exp(-a_2 \log p)$,*

$$\|\hat{\beta} - \beta^*\|_2 \leq a_1^* \sqrt{s \log p/n} \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_1 \leq 4a_1^* s \sqrt{\log p/n}.$$

- (ii) *(Soft sparsity case) Assume β^* satisfies the soft sparsity assumption (4). For any R satisfying $R \geq \sqrt{\log p/n}$ and $a_1^* \sqrt{\log p/n} \max\{2, R\} < 1/2$, we have, with*

probability at least $1 - 4\delta_n - 4\exp(-\log p) - 2a_1 \exp(-a_2 \log p)$,

$$\begin{aligned} \|\widehat{\beta} - \beta^*\|_2 &\leq a_2^* R^{1/2} (\log p/n)^{1/4} \quad \text{and} \\ \|\widehat{\beta} - \beta^*\|_1 &\leq 4(a_2^* \sqrt{|S_a|} R^{1/2} (\log p/n)^{1/4} + \|\beta_{S_a^c}^*\|_1), \end{aligned}$$

where $a_2^* = 2 \max \{\sqrt{2}a_1^*, \sqrt{a_1^*}\}$, $S_a = S_{-a} \cup \{1\}$ with $S_{-a} = \{j : |\beta_j^*| > a, 2 \leq j \leq p\}$, and $a > 0$ is an arbitrary thresholding parameter.

Remark 3 (On the results of QR-LASSO for the hard sparsity case). In this case, the L_2 estimation error of QR-LASSO has the rate $\sqrt{s \log p/n}$. This matches the minimax optimal rate for LS-LASSO, established in Raskutti, Wainwright, and Yu (2011) under the assumption of sub-Gaussian errors for the hard sparsity case. In the oracle case (when the underlying model is known), the L_2 estimation error has the rate $\sqrt{s/n}$. The above minimax rate is near-oracle up to a factor of order $\sqrt{\log p}$, the price to pay for not knowing in advance which of the p covariates are relevant.

The results in Theorem 2(i) for the hard sparsity case are inspired by the earlier work of Belloni and Chernozhukov (2011) and Wang (2013), which obtained the same rates for the L_2 error bound. However, our proof employs a different technique and requires weaker conditions. Comparing with the conditions in Belloni and Chernozhukov (2011), we relaxed the conditions on both Σ and ϵ_i . We have dropped their *restricted nonlinearity condition* on Σ (their condition D.4), which would require $q := \inf_{\delta \in A, \delta \neq 0} \frac{\{E(|X_i^T \delta|^2)\}^{3/2}}{E(|X_i^T \delta|^3)} > 0$ for some restricted set A . Such a condition is not needed for the parallel theory of LS-LASSO. Furthermore, if the nonlinear impact coefficient q converges to zero at a sufficiently fast rate, this may have a negative impact on the feasible range of n and p through the *growth condition* $\sqrt{s \log(p \vee n)} \leq O(q\sqrt{n})$ required in the main theorem (Theorem 2) of Belloni and Chernozhukov (2011). Unlike Belloni and Chernozhukov (2011), we do not require the conditional random error density $f_i(t)$ to be continuously differentiable nor the derivative to be uniformly bounded everywhere. We only need a uniform lower bound for $f_i(t)$ in a neighborhood of zero. Our assumptions are also significantly weaker than those in Wang (2013), which required independent and identically distributed (i.i.d.) random errors and a restricted isometry-type condition in addition to the restricted eigenvalue condition.

Remark 4 (On the results of QR-LASSO for the soft sparsity case). The results in Theorem 2(ii) for the soft sparsity case are new for high-dimensional QR. The soft sparsity scenario allows for dense small coefficients. The radius of the L_1 -ball $\mathbb{B}_1(R)$ is allowed to shrink or diverge with the sample size n . In this case, we obtain the L_2 error rate $R^{1/2} (\log p/n)^{1/4}$ for QR-LASSO, which also matches the minimax optimal rate in the soft sparsity case for LS-LASSO in Raskutti et al. (2011). The L_1 -error bound is larger than the L_2 error bound. However, one may still achieve an L_1 -consistency rate under additional structural assumptions on β^* . As suggested by an anonymous referee, let us consider an example of an approximately sparse

model for illustration. Without loss of generality, we assume $|\beta_2^*| \geq |\beta_3^*| \geq \dots \geq |\beta_p^*|$, and that there exists a positive integer $q < p$ such that $q\sqrt{\log p/n} = o(1)$ and $\beta_j = (\frac{1}{2})^{j-q}\sqrt{\log p/n}$, for $j = q+1, \dots, p$. Then $\sum_{j=q+1}^p |\beta_j^*| = O(\sqrt{\log p/n})$. Taking $a = |\beta_q^*|$ and assuming R is bounded, the result in Theorem 2(ii) implies that $\|\hat{\beta} - \beta^*\|_1 = O(q^{1/2}R^{1/2}(\log p/n)^{1/4})$.

Remark 5. The regularization parameter λ is taken to be of the order $\sqrt{\log p/n}$, the universal penalty level introduced in Donoho and Johnstone (1994). The literature of regularized high-dimensional regression often focuses on statistical analysis with a penalty parameter of this order (e.g., Bickel et al., 2009). In practice, an appealing approach (Belloni and Chernozhukov, 2011) is to directly simulate λ as the $(1 - \alpha)$ -quantile of the distribution of $c\|\tilde{S}_n\|_\infty$, for some small $\alpha > 0$. This is feasible by observing that the distribution of $\|\tilde{S}_n\|_\infty$ is pivotal. With this simulated λ , the same estimation error bound would hold with probability at least $1 - \alpha - 4\delta_n - 4\exp(-\log p) - 2a_1\exp(-a_2\log p)$.

Remark 6. Lemma B.1 in Appendix B demonstrates that for a certain class of heavy-tailed error distributions and a general class of design matrices, there is a positive probability that LS-LASSO cannot achieve the near-oracle rate derived under the normality assumption no matter the choice of the tuning parameter. In contrast, the results for regularized QR in this paper hold for a much larger class of error distributions. Our results, hence, provide strong evidence for the robustness and broader applicability of QR in high dimensions.

Theorem 3 gives the estimation error bounds for the feasible local solutions of QR-NCP. The nonconvex penalty function is assumed to satisfy the following general conditions. The penalty function $p_\lambda(t)$ is defined on the real line and is symmetric about zero. It is assumed to be nondecreasing and concave for $t \in [0, +\infty)$, with a continuous derivative $p'_\lambda(t)$ on $(0, +\infty)$ and $\lim_{t \rightarrow 0+} p'_\lambda(t) = \lambda$. For $t > 0$, $p_\lambda(t)$ is nonincreasing in t . Furthermore, there exists a constant $\gamma_0 > 0$ such that the function $t \mapsto p_\lambda(t) + \frac{\gamma_0}{2}t^2$ is convex. This class of nonconvex penalty functions, in particular, includes the popular choices SCAD and MCP penalties discussed in Section 2.1.

THEOREM 3 (QR-NCP). *Let $\lambda = k_0\sqrt{\log p/n}$, where $k_0 \geq 4\max\{2\sqrt{m_x}, c^*\}$. Suppose conditions (C1)–(C3) are satisfied. Consider any feasible local solution $\hat{\beta}$ such that $\|\hat{\beta}\|_1 < \kappa$, for some $\kappa > \|\beta^*\|_1$, and the KKT condition (7) is satisfied.*

- (i) *(Hard sparsity case) Assume β^* satisfies the hard sparsity assumption (3). If $\sqrt{\log p/n} < \frac{2a^*}{3\kappa k_0}$ and $a^* > \frac{3}{4}\gamma_0$, then, with probability at least $1 - 4\delta_n - 4\exp(-\log p) - 2a_1\exp(-a_2\log p)$,*

$$\|\hat{\beta} - \beta^*\|_2 \leq a_3^* \sqrt{s \log p/n}, \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_1 \leq 4a_3^* s \sqrt{\log p/n},$$

$$\text{where } a_3^* = \frac{6k_0}{4a^* - 3\gamma_0}.$$

- (ii) *(Soft sparsity case) Assume β^* satisfies the soft sparsity assumption (4). If $\sqrt{\log p/n} < \max\{R, \frac{2a^*}{3\kappa k_0}\}$ and $a^* > \gamma_0$, then, with probability at least $1 - 4\delta_n -$*

$$4\exp(-\log p) - 2a_1 \exp(-a_2 \log p),$$

$$\|\widehat{\beta} - \beta^*\|_2 \leq a_4^* R^{1/2} (\log p/n)^{1/4}, \quad \text{and}$$

$$\|\widehat{\beta} - \beta^*\|_1 \leq 4a_4^* \sqrt{|S_a|} R^{1/2} (\log p/n)^{1/4} + 2\|\beta_{S_a^c}^*\|_1,$$

where $a_4^* = 2 \max \left\{ \frac{3\sqrt{2}k_0}{2(a^* - \gamma_0)}, \sqrt{\frac{k_0}{a^* - \gamma_0}} \right\}$, and $|S_a|$ is defined the same as in Theorem 2.

Remark 7. The assumption of the KKT condition (7) is satisfied for the local solutions corresponding to the SCAD or MCP penalty function. The side condition $\|\widehat{\beta}\|_1 \leq \kappa$ as was also adopted in Loh and Wainwright (2015) in order to focus on sensible local solutions. The high-probability error bound in this theorem applies to any feasible local solution within the radius κ of the true value β^* . A careful examination of the proof reveals that κ can diverge to ∞ as long as $\kappa \sqrt{\log p/n} = o(1)$. It is noted that we do not restrict the local solution to be sparse even for the hard sparsity setting. In practice, one would not wish to choose an exceedingly large κ , as the conditions of the theorem suggest that the non-asymptotic bounds hold with a larger n when κ is larger.

Embedded in the proof of Theorem 3 is a result on the geometric structure of the local solution of QR-NCP. Unlike the result in Lemma 2 about the geometric structure of the global solution of QR-LASSO, the result for QR-NCP in Corollary 1 is new.

COROLLARY 1 (QR-NCP). *Assume the conditions of Theorem 3 are satisfied. Then, with probability at least $1 - 4\delta_n - 4\exp(-\log p) - 2a_1 \exp(-a_2 \log p)$, where a_1 and a_2 are the positive constants in Theorem 1, we have (i) $\widehat{v} \in \widetilde{\Gamma}_H$ under the hard sparsity assumption and (ii) $\widehat{v} \in \widetilde{\Gamma}_W$ under the soft sparsity assumption.*

4. QUANTILE PREDICTION ERROR BOUNDS

Finally, we establish theoretical guarantees for the prediction error bounds of the regularized QR estimators in high dimensions. The empirical quantile prediction error is given by

$$R_n(\widehat{\beta}) = Q_n(\widehat{\beta}) - Q_n(\beta^*),$$

where the sample quantile loss function $Q_n(\beta)$ is defined in (5). Our result does not impose restrictions on which regularized procedure is used. Specifically, the prediction error is evaluated using $\widehat{\beta}$, either the global solution from QR-LASSO or a local solution from QR-NCP.

A key result of this section is a general characterization of the prediction error bound based on subgradients. Consider any two subgradients S_n and \bar{S}_n in $\partial Q_n(\beta)$. The convexity of the loss function $Q_n(\beta)$ implies that

$$Q_n(\widehat{\beta}) \geq Q_n(\beta^*) + S_n(\beta^*)^T (\widehat{\beta} - \beta^*), \quad (17)$$

$$Q_n(\beta^*) \geq Q_n(\hat{\beta}) + \bar{S}_n(\hat{\beta})^T(\hat{\beta} - \beta^*). \quad (18)$$

Combining inequalities (17) and (18), we immediately obtain

$$|Q_n(\hat{\beta}) - Q_n(\beta^*)| \leq \max\{\|S_n(\beta^*)\|_\infty, \|\bar{S}_n(\hat{\beta})\|_\infty\} \|\hat{\beta} - \beta^*\|_1.$$

This leads to simple bounds for the prediction error for the general case (global or local solution, hard or soft sparsity).

THEOREM 4. *Let $\lambda = k_0 \sqrt{\log p/n}$, where $k_0 \geq 4 \max\{2\sqrt{m_x}, c^*\}$. With probability at least $1 - \delta_n - 2 \exp(-\log p)$, we have*

$$|R_n(\hat{\beta})| \leq 4\lambda \|\hat{\beta} - \beta^*\|_1. \quad (19)$$

Corollary 2 summarizes the results for the quantile prediction error. These results are new for high-dimensional QR. Our results extend the prediction error bound for the LS-LASSO (see Greenshtein and Ritov, 2004; Bunea, Tsybakov, and Wegkamp, 2007; Bickel et al., 2009; Raskutti et al., 2011). Unlike QR, the precision error bound for LS-LASSO is usually studied based on the least-squares loss function $n^{-1} \sum_{i=1}^n (\epsilon_i - X_i^T \gamma)^2$. Our approach can also be applied to other (possibly nondifferentiable) convex loss functions.

COROLLARY 2.

(i) *(Slow rate, without sparsity assumption) Let $\hat{\beta}$ be the QR-LASSO or QR-NCP estimator. For any tuning parameter λ , for any n , we have*

$$\begin{aligned} |R_n(\hat{\beta})| &\leq 2\lambda \|\beta^*\|_1, \quad \text{for LASSO,} \\ |R_n(\hat{\beta})| &\leq 4\lambda (\|\beta^*\|_1 + R), \quad \text{for QR-NCP.} \end{aligned}$$

(ii) *(QR-LASSO, faster rate) Let $\hat{\beta}$ be the QR-LASSO estimator. Assume the conditions of Theorem 2 are satisfied. Then, with probability at least $1 - 4\delta_n - 4 \exp(-\log p) - 2a_1 \exp(-a_2 \log p)$,*

$$\begin{aligned} |R_n(\hat{\beta})| &\leq b_1^* s \log p/n, \quad \text{hard sparsity case,} \\ |R_n(\hat{\beta})| &\leq b_2^* (a_2^* \sqrt{|S_a|} R^{1/2} (\log p/n)^{3/4} + \|\beta_{S_a}^*\|_1 (\log p/n)^{1/2}), \\ &\quad \text{soft sparsity case,} \end{aligned}$$

where $b_1^* = 16k_0 a_1^*$, $b_2^* = 16k_0$, and a_1^* , a_2^* , and S_a are defined the same as in Theorem 2.

(iii) *(QR-NCP, faster rate) Let $\hat{\beta}$ be the QR-NCP estimator. Assume the conditions of Theorem 3 are satisfied. Then, with probability at least $1 - 4\delta_n - 4 \exp(-\log p) - 2a_1 \exp(-a_2 \log p)$,*

$$|R_n(\widehat{\beta})| \leq b_3^* s \log p/n, \quad \text{hard sparsity case,}$$

$$|R_n(\widehat{\beta})| \leq b_4^* [2a_4^* \sqrt{|S_a|} R^{1/2} (\log p/n)^{3/4} + \|\beta_{S_a^c}^*\|_1 (\log p/n)^{1/2}],$$

soft sparsity case,

where $b_3^* = 4k_0 a_3^*$, $b_4^* = 8k_0$, and a_3^* , a_4^* , and S_a are defined the same as in Theorem 3.

Remark 8. The rate in (i) is known as the “slow rate” for prediction error in the literature on LS-LASSO. It is obtained without assuming any structure for β^* or any assumption on the design matrix. In (ii) and (iii), the rates for the hard sparsity case are the same as the so-called “fast rate” in the literature for LS-LASSO. The upper bounds for the soft sparsity case in (ii) and (iii) permit faster rates when the true value β^* has certain desirable structural properties, particularly when the number of relatively large signals in β^* is of a relatively small order, whereas the number of relatively small signals is much smaller relative to R . For example, for the example of the approximately sparse model discussed in Remark 4, we have $|R_n(\widehat{\beta})| = O(q^{1/2} R^{1/2} (\log p/n)^{3/4})$ in the soft sparsity case for both QR-LASSO and QR-NCP. Finally, we note that for all these scenarios, we require overall milder conditions on the random error distributions than those required in the literature for LS-LASSO.

5. A MONTE CARLO EXPERIMENT

In this section, we carry out a Monte Carlo experiment to confirm some of the theoretical findings about LASSO or SCAD regularized QR versus least-squares regression under the same form of penalization.

We generate (X_1, X_2, \dots, X_p) from the multivariate normal distribution $N_p(0, \Sigma)$ with $\Sigma = (\sigma_{jk})_{p \times p}$ and $\sigma_{jk} = 0.5^{|j-k|}$, $1 \leq j, k \leq p$. For the regression parameter β^* , we consider two different models.

- Model 1 (sparser model): $\beta^* = (2, 1, 1.5, 1.75, 0_{p-4}^T)^T$,
- Model 2 (denser model): $\beta^* = \frac{3}{n} (1_n, 0_{p-n}^T)^T$,
- Model 3 (dense model): $\beta^* = (2, 1, 1.5, 1.75, \beta_5^*, \dots, \beta_p^*)^T$, where $\beta_j^* = (0.5)^{j-4} \sqrt{\log p/n}$, for $j = 5, \dots, p$,

where 0_k denotes the k -dimensional vector of zeros, whereas 1_k denotes the k -dimensional vector of ones. For each model, we consider three different random error distributions for ϵ_i : the $N(0, 1)$ distribution, the mixture normal distribution $aN(0, 1) + (1-a)N(0, 10^2)$ where $a \sim \text{Bernoulli}(0.95)$, and the Cauchy distribution with location 0 and scale 1.

We consider LS_Oracle, QR_Oracle, LS_LASSO, QR_LASSO, LS_SCAD (SCAD regularized least-squares regression), and QR_SCAD (SCAD regularized QR), where the quantile methods are based on $\tau = 0.5$, and LS_Oracle and QR_Oracle are computed using the true model structure. For $n = 100$ and $p = 500$,

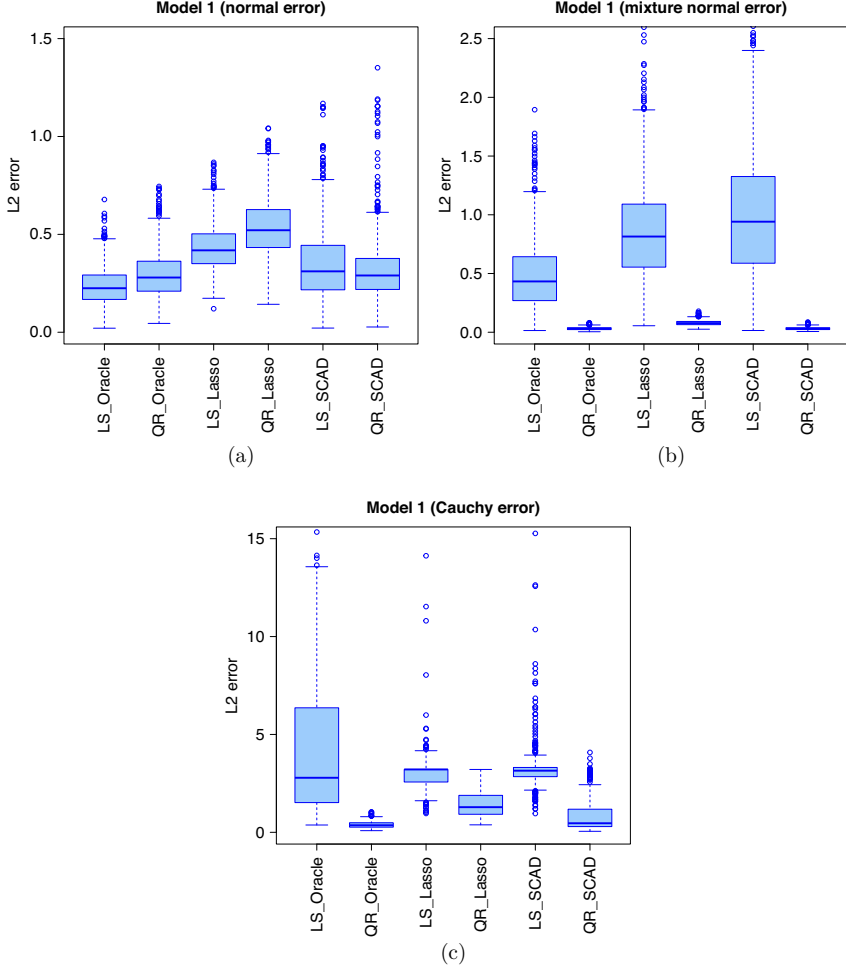


FIGURE 1. Box plots of the L_2 estimation error for Model 1 based on six methods: LS_Oracle, QR_Oracle, LS_LASSO, QR_LASSO, LS_SCAD, and QR_SCAD.

the box plots for the L_2 -errors of the six methods based on 1,000 simulation runs for Models 1–3 are given in Figures 1–3, respectively.

For the LASSO penalty, the tuning parameter is selected using a fivefold cross-validation. For the SCAD penalty, the tuning parameter is selected using a high-dimensional BIC procedure (Wang, Kim, and Li, 2013b; Lee, Noh, and Park, 2014, among others). We observe from Figure 1 that for normal random errors, LS-LASSO is slightly more efficient than QR-LASSO, but its performance deteriorates substantially for the mixture normal random errors and Cauchy errors. The nonconvex SCAD penalty leads to smaller L_2 error than the LASSO penalty.

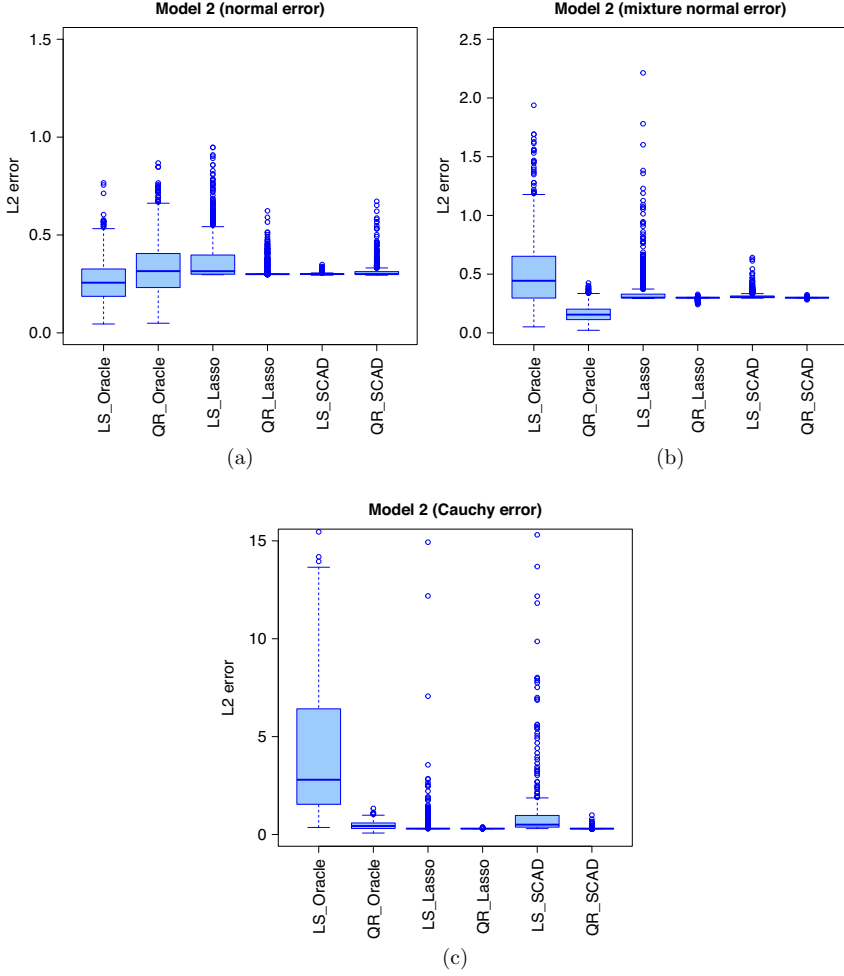


FIGURE 2. Box plots of the L_2 estimation error for Model 2 based on six methods: LS_Oracle, QR_Oracle, LS_LASSO, QR_LASSO, LS_SCAD, and QR_SCAD.

Figure 2 suggests that for the denser model under consideration, QR-LASSO and QR-SCAD have similar performance to that of LS-LASSO and LS-SCAD for the normal errors and have much smaller error for the mixture normal errors and Cauchy errors. Similar observations are obtained from Figure 3 for the dense model.

Finally, to investigate the effects of different sample sizes, we consider Model 1 with the mixture normal error distribution as the sample size varies between 100 and 800. Figure 4 depicts the L_2 errors for LS_LASSO, QR_LASSO, LS_SCAD, and QR_SCAD in this setting. The plot suggests that the L_2 error decreases as n

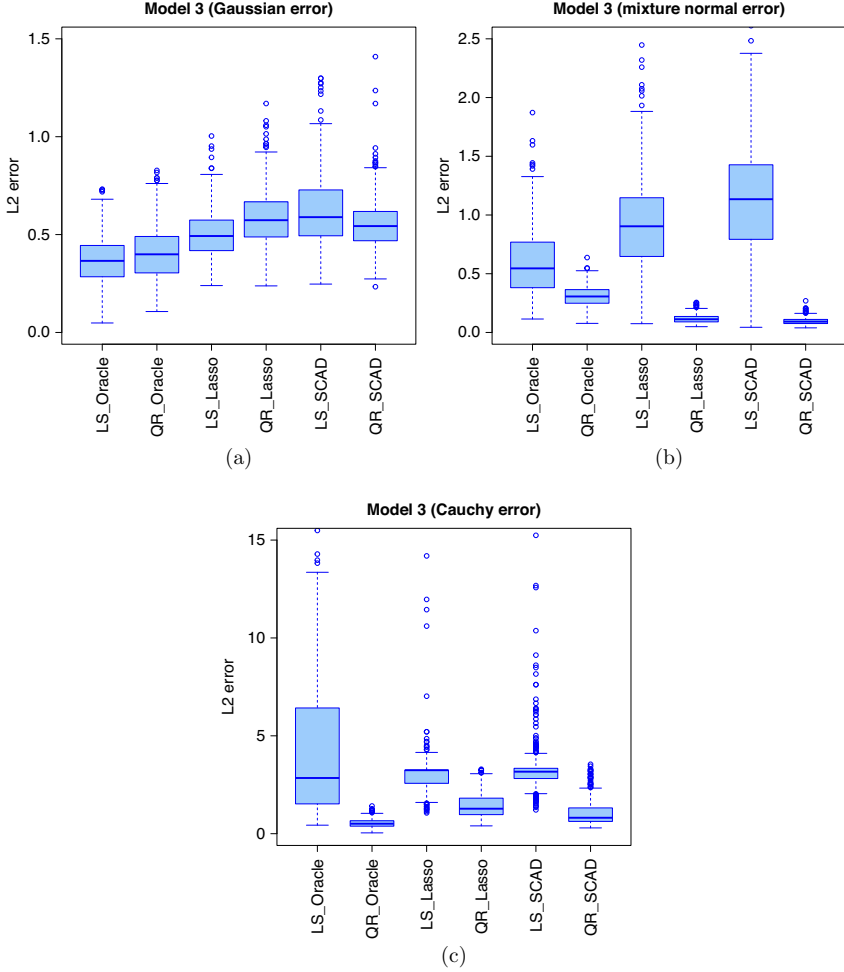


FIGURE 3. Box plots of the L_2 estimation error for Model 3 based on six methods: LS_Oracle, QR_Oracle, LS_LASSO, QR_LASSO, LS_SCAD, and QR_SCAD.

increases. It also suggests that for the heavy-tailed error distribution, regularized QR substantially outperforms regularized least-squares regression.

6. DISCUSSION

By developing a new and unified subgradient approach, we present several significant results on the fundamental properties of regularized QR in high dimensions, where the number of covariates can grow at an exponential rate of the sample size. We demonstrate that QR-LASSO enjoys the near-oracle rate in estimation error for

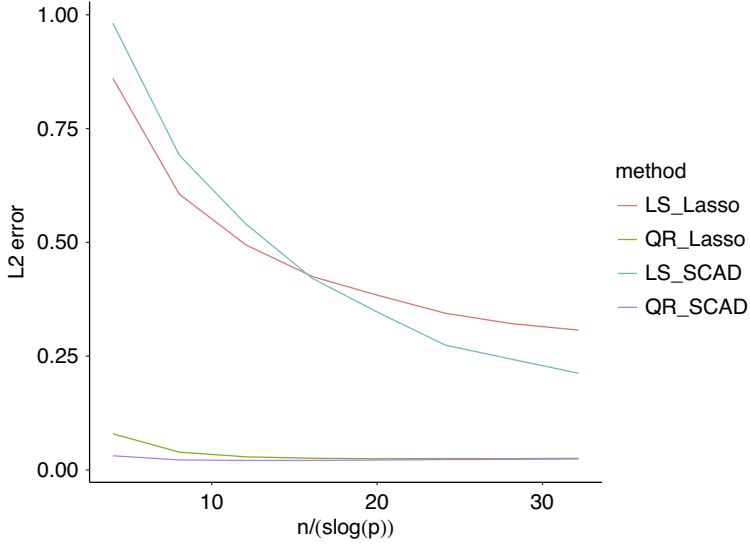


FIGURE 4. The L_2 estimation error with varying sample sizes (Model 1, mixture normal error).

a much richer class of error distributions than LS-LASSO does. The result renders theoretical support for the wide applicability of QR in high-dimensional problems. Our analysis is carried out under both the hard sparsity and the soft sparsity models. We also prove that any feasible local solution of QR with some commonly used nonconvex penalty functions, such as SCAD or MCP, enjoys the same estimation error rates, making the theoretical results more meaningful to practical algorithms which often yield a local optimum instead of a global one.

The unified approach based on subgradients can be useful for studying other high-dimensional problems involving a nonsmooth loss function. Hence, it is of independent interest. For example, the support vector machine (SVM) is a powerful binary classification tool with high accuracy and great flexibility. The sample loss function for SVM is given by

$$n^{-1} \sum_{i=1}^n (1 - Y_i \mathbf{x}_i^T \boldsymbol{\beta})_+,$$

where $Y_i \in \{1, -1\}$, and $(1 - u)_+ = \max\{1 - u, 0\}$ denotes the hinge loss. As another example, we consider the rank loss function for robust high-dimensional regression (e.g., Wang et al., 2020) given by

$$[n(n-1)]^{-1} \sum_{i \neq j} |(Y_i - \mathbf{x}_i^T \boldsymbol{\beta}) - (Y_j - \mathbf{x}_j^T \boldsymbol{\beta})|.$$

Both examples involve nonsmooth loss functions. The subgradient techniques developed in this paper can be applied to study the statistical properties of the global and local solutions of the corresponding regularized estimation problems in high dimensions when $p \gg n$.

The results of the paper are also useful for inference procedures for high-dimensional QR, which need an initial value (see, for example, Belloni, Chernozhukov, and Kato, 2014; Zhao, Kolar, and Liu, 2014; Bradic and Kolar, 2017; Belloni, Chernozhukov, and Kato, 2019).

APPENDIX A. Parameter Identification in High Dimensions

Let $\mathbb{X} = (\mathbb{X}_1, \dots, \mathbb{X}_p)$ be the $n \times p$ matrix of covariates, where \mathbb{X}_j denotes the j th ($j = 1, \dots, p$) column of the matrix of covariates. In the classical asymptotic framework where p is smaller than n , it is usually assumed that \mathbb{X} has a full column rank, which allows to uniquely identify β^* . In contrast, in the high-dimensional scenario where $p \gg n$, β^* is generally not identifiable in the absence of additional structural assumptions as \mathbb{X} has at most column rank n .

In high dimensions, β^* in general is not uniquely defined. Suppose model (1) is satisfied by $\beta^* = \beta_0^*$. Consider the affine space $\{\beta^* \in \mathbb{R}^p : \mathbb{X}\beta^* = \mathbb{X}\beta_0^*\}$. We emphasize that the error bounds derived in this paper apply to any β^* from this affine space and does not require the unique identification.

Let $\text{Ker}(\mathbb{X}) = \{\beta \in \mathbb{R}^p : \mathbb{X}\beta = 0\}$ be the null space of \mathbb{X} . If β^* satisfies (1), then $\beta^* + \beta$ also satisfies (1), $\forall \beta \in \text{Ker}(\mathbb{X})$. The extent of identifiability can be measured by the diameter of the set $N_0(\mathbb{X}) = \text{Ker}(\mathbb{X}) \cap \mathbb{B}$, defined as $\max_{\beta \in N_0(\mathbb{X})} \|\beta\|_2$, where $\mathbb{B} = \mathbb{B}_0(s)$ under the hard sparsity assumption, whereas $\mathbb{B} = \mathbb{B}_1(R)$ under the soft sparsity assumption. The following lemma characterizes the properties of the diameter of $N_0(\mathbb{X})$.

LEMMA A.1. Assume the vector $X_{-i} = (X_{i2}, \dots, X_{ip})'$ is a mean-zero sub-Gaussian random vector.

- (i) (Hard sparsity case) Assume $\eta_{\min}(s) = \inf_{v: \|v\|_2=1, \|v\|_0 \leq s} v' \Sigma v > 0$, where $\Sigma = E(X_i X_i')$. Then

$$P\left(\max_{\beta \in N_0(\mathbb{X})} \|\beta\|_2 = 0\right) \geq 1 - \alpha_1^* \exp(-\alpha_2^* \log p),$$

where α_1^* and α_2^* are positive constants.

- (ii) (Soft sparsity case) Assume $\log(p) = o(n)$. Then

$$P\left(\max_{\beta \in N_0(\mathbb{X})} \|\beta\|_2 \leq \alpha' \xi_{\min}^{-1/2} R \sqrt{\log p / n}\right) \geq 1 - \alpha_1^* \exp(-\alpha_2^* \log p),$$

for all n sufficiently large, where α_1^* , α_2^* , and α' are positive constants, and ξ_{\min} is the smallest eigenvalue of Σ .

The above lemma can be considered as a generalization of Lemma 1 of Raskutti et al. (2011) to the random covariates case. From this lemma, it can be seen that under some general conditions, $\text{Ker}(\mathbb{X}) = \{0\}$ with high probability for the hard sparsity case, i.e., the

sparse β^* satisfying (1) is unique, whereas for the soft sparsity case, $\text{Ker}(\mathbb{X})$ is a shrinking neighborhood around 0 with high probability if $\xi_{\min}^{-1/2} R\sqrt{\log p/n} \rightarrow 0$.

Proof of Lemma A.1. (i) (Hard sparsity case) Consider $\mathbb{C}_1(s) = \{\theta \in \mathbb{R}^p : \|\theta\|_2 = 1, \|\theta\|_0 \leq 2s\}$. Applying Lemma D.1 with $t = \eta_{\min}(s)/2$, we have

$$\begin{aligned} & P\left\{ \sup_{\theta \in \mathbb{C}_1(s)} |n^{-1} \|\mathbb{X}\theta\|_2^2 - E(n^{-1} \|\mathbb{X}\theta\|_2^2)| \geq \eta_{\min}(s)/2 \right\} \\ & \leq \alpha_2 \exp\left(-\alpha_1 n \min(\eta_{\min}^2(s)/(4\sigma_x^4), \eta_{\min}(s)/(2\sigma_x^2)) + 2s \log p\right), \end{aligned} \quad (\text{A.1})$$

for some positive constants α_1 and α_2 . We argue by contradiction and assume $\max_{\beta \in N_0(\mathbb{X})} \|\beta\|_2 \neq 0$. Then there exists a $\theta \neq 0$ in $\mathbb{B}_0(s)$ such that $\mathbb{X}\theta = 0$. Let $\tilde{\theta} = \theta/\|\theta\|_2$, then $\tilde{\theta} \in \mathbb{C}_1(s)$. It follows from (A.1) that there exist some positive constants α_1^* and α_2^* such that with probability at least $1 - \alpha_1^* \exp(-\alpha_2^* n)$,

$$n^{-1} \|\mathbb{X}\tilde{\theta}\|_2^2 \geq \tilde{\theta}' \Sigma \tilde{\theta} - \eta_{\min}(s)/2 \geq \eta_{\min}(s)/2 > 0,$$

which contradicts the assumption $\mathbb{X}\tilde{\theta} = 0$.

(ii) (Soft sparsity case) By Lemma D.1, there exist some positive constants α_1 and α_2 such that $\forall t > 0$,

$$\left\{ \sup_{\theta \in \mathbb{C}(k)} |n^{-1} \|\mathbb{X}\theta\|_2^2 - E(n^{-1} \|\mathbb{X}\theta\|_2^2)| \geq t \right\} \leq \alpha_2 \exp(-\alpha_1 n \min(t^2/\sigma_x^4, t/\sigma_x^2) + 2k \log p),$$

where $\mathbb{C}(k) = \{\theta \in \mathbb{R}^p : \|\theta\|_2 \leq 1, \|\theta\|_0 \leq 2k\}$ for an arbitrary $k \geq 1$. Taking $t = \xi_{\min}/54$ and $k = \frac{1}{4}\alpha_1 \min(t^2/\sigma_x^4, t/\sigma_x^2) \frac{n}{\log p}$, for all n sufficiently large, we have

$$P\left\{ \sup_{\theta \in \mathbb{C}(k)} |n^{-1} \|\mathbb{X}\theta\|_2^2 - E(n^{-1} \|\mathbb{X}\theta\|_2^2)| < t \right\} \geq 1 - \alpha_2 \exp(-\alpha_1 n \min(t^2/\sigma_x^4, t/\sigma_x^2)/2).$$

Then the same argument as in Lemmas 12 and 13 of Loh and Wainwright (2012) for their (F.8) leads to

$$n^{-1} \|\mathbb{X}\theta\|_2^2 \geq \frac{\xi_{\min}}{2} \|\theta\|_2^2 - c' n^{-1} (\log p) \|\theta\|_1^2,$$

for all $\theta \in \mathbb{R}^p$, with probability at least $1 - \alpha_1^* \exp(-\alpha_2^* n)$, for all n sufficiently large, for some positive constants c' , α_1^* , and α_2^* . This implies that, for any $\theta \in \text{Ker}(\mathbb{X}) \cap \mathbb{B}(R)$,

$$0 = n^{-1} \|\mathbb{X}\theta\|_2^2 \geq \frac{\xi_{\min}}{2} \|\theta\|_2^2 - c' R^2 n^{-1} \log p,$$

or $\max_{\theta \in N_0(\mathbb{X})} \|\theta\|_2 \leq \alpha' \xi_{\min}^{-1/2} R\sqrt{\log p/n}$ for some positive constant α' with probability at least $1 - \alpha_1^* \exp(-\alpha_2^* n)$, for all n sufficiently large. Hence, the conclusion follows. ■

APPENDIX B. Performance Lower Bound for LASSO with Heavy-Tailed Errors

The LS-LASSO estimator $\hat{\beta}^{LS}$ is defined as

$$\hat{\beta}^{LS} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ (2n)^{-1} \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \|\beta_-\|_1 \right\}. \quad (\text{B.1})$$

Suppose $\epsilon_i = Y_i - X_i^T \beta^*$ ($i = 1, \dots, n$) are independent Cauchy(0,1) random variables with the density function $f(\epsilon) = \frac{1}{\pi(1+\epsilon^2)}$, and the characteristic function $\phi(t) = E[e^{it\epsilon}] = e^{-|t|}$, $\forall t \in \mathbb{R}$.

Assume $\min_{1 \leq j \leq p} E(|x_{ij}|) \geq \zeta_* > 0$. The $\hat{\beta}^{LS}$ estimator is called a nondegenerate solution if it has at least one nonzero component. We consider the asymptotic regime where $\log p = o(n)$.

LEMMA B.1. *Consider the setting described above. Let a_0 be any positive constant, and consider an arbitrary $\lambda \in (0, a_0)$. Let $\hat{\beta}^{LS}$ be any nondegenerate solution of (B.1) corresponding to λ . Then there exist some constants $\zeta_1 > 0$ and $0 < \zeta_2 < 1$ such that if (n, p) satisfies $p \geq 5$, $\log p \leq n/4$, and $\sqrt{\log p/n} \leq \zeta_1$, then*

$$P(\|\hat{\beta}^{LS} - \beta^*\|_1 > 1) \geq \zeta_2,$$

where the constants ζ_1 and ζ_2 do not depend on (n, p) .

Loh (2017) showed that if one chooses the tuning parameter λ at the regular rate $\sqrt{\log p/n}$, which is the theoretical choice that leads to the near-oracle performance of LS-LASSO for normal random errors, then the KKT condition for LS-LASSO may not hold for Cauchy random errors. The above lemma strengthens the result by showing that a different choice of λ cannot fix the problem. Indeed, for any $\lambda \in (0, a_0)$, the L_1 -estimation error of LS-LASSO has a positive probability to exceed one. This paper shows that QR-LASSO is still consistent in this case with $\lambda \sim \sqrt{\log p/n}$. This lemma, hence, renders strong support for the robustness of QR-LASSO in high dimensions. Related to this, Fan et al. (2014) showed that for a specially designed fixed design matrix, LS-LASSO cannot be sign-consistent unless a certain signal condition is satisfied. Fan, Li, and Wang (2017) investigated the estimation of high-dimensional mean regression in the absence of symmetry and light-tailed assumptions, but their conditions exclude Cauchy random errors.

It is worth noting that the above probability bound is nonasymptotic and holds for all n sufficiently large and does not become smaller as n increases. Although this lemma focuses on the Cauchy random error for a clean presentation, analogous inconsistency results hold more generally for the class of α -stable distributions with $\alpha \in (0, 2)$. Specifically, ϵ_i has an α -stable distribution with scale parameter ξ if its characteristic function $E\{\exp(it\epsilon_i)\} = \exp(-\xi^\alpha |t|^\alpha)$, $\forall t \in \mathbb{R}^p$. The standard Cauchy distribution is an α -stable distribution with $\alpha = \xi = 1$ (see Nolan (2003) for an introduction to stable distributions).

Proof of Lemma B.1. Assume the contrary is true, that is, $\|\hat{\beta}^{LS} - \beta^*\|_1 \leq 1$ for some $\lambda \in (0, a_0)$. As $\hat{\beta}^{LS}$ is a nondegenerate solution, it has at least one nonzero component. Without loss of generality, we assume $\hat{\beta}_j^{LS} \neq 0$, for some $1 \leq j \leq p$. By the KKT condition, $\hat{\beta}^{LS}$ must satisfy

$$e_j^T (n^{-1} \mathbb{X}^T \mathbb{X}) (\beta^* - \hat{\beta}^{LS}) + e_j^T n^{-1} \mathbb{X}^T \epsilon + \lambda \text{sign}(\hat{\beta}_j^{LS}) = 0, \quad (\text{B.2})$$

where e_j is a p -dimensional unit vector with the j entry being one and all the other entries being zero.

Consider the event $\Omega_{n1} = \{\|\hat{\Sigma}\|_\infty \leq 12\zeta_0^2\}$. By Lemma D.2, $P(\Omega_{n1}) \geq 1 - 2\exp(-\log p)$. On Ω_{n1} ,

$$\begin{aligned}
|e_j^T(n^{-1}\mathbb{X}'\mathbb{X})(\beta^* - \widehat{\beta}^{LS})| &\leq \|e_j^T(n^{-1}\mathbb{X}'\mathbb{X})\|_\infty \|(\beta^* - \widehat{\beta}^{LS})\|_1 \\
&\leq \|\widehat{\Sigma}\|_\infty \|(\beta^* - \widehat{\beta}^{LS})\|_1 \\
&\leq 12\zeta_0^2 \|(\beta^* - \widehat{\beta}^{LS})\|_1.
\end{aligned}$$

As $|\lambda \text{sign}(\widehat{\beta}_j^{LS})| \leq a_0$, it follows from the KKT condition (B.2) that on the event Ω_{n1} , we have $|e_j^T n^{-1}\mathbb{X}^T \epsilon| \leq a_0 + 12\zeta_0^2$.

Conditional on \mathbb{X}_j , $e_j^T n^{-1}\mathbb{X}^T \epsilon = n^{-1}\mathbb{X}_j^T \epsilon$ has a Cauchy(0, $n^{-1} \sum_{i=1}^n |x_{ij}|$) distribution, by checking the form of its characteristic function. By the property of Cauchy distribution, we have

$$P(|e_j^T n^{-1}\mathbb{X}^T \epsilon| > a_0 + 12\zeta_0) = 2E_{\mathbb{X}_j} \left\{ \frac{1}{2} - \frac{1}{\pi} \arctan \left(\frac{a_0 + 12\zeta_0}{n^{-1} \sum_{i=1}^n |x_{ij}|} \right) \right\}.$$

Note that $n^{-1} \sum_{i=1}^n |x_{ij}|$ on Ω_n , $e_j^T n^{-1}\mathbb{X}^T \epsilon$ has a Cauchy(0, b) distribution with $0 < b < 12\zeta_0$. Consider the event $\Omega_{n2} = \{ \min_{1 \leq j \leq p} n^{-1} \sum_{i=1}^n |x_{ij}| \geq \zeta_*/2 \}$. By Lemma D.3, $P(\Omega_{n2}) \geq 1 - 2 \exp(-\log p)$. On the event Ω_{n2} ,

$$P(|e_j^T n^{-1}\mathbb{X}^T \epsilon| > a_0 + 12\zeta_0) \geq 1 - \frac{2}{\pi} \arctan(\zeta_*^{-1}(2a_0 + 24\zeta_0)).$$

Hence, on the event $\Omega_{n1} \cap \Omega_{n2}$, with probability at least $1 - \frac{2}{\pi} \arctan(\zeta_*^{-1}(2a_0 + 24\zeta_0))$, there is a contradiction. We thus have

$$\begin{aligned}
P(\|\beta^* - \widehat{\beta}^{LS}\|_1 > 1) &\geq P(\|\beta^* - \widehat{\beta}^{LS}\|_1 > 1 | \Omega_{n1} \cap \Omega_{n2}) P(\Omega_{n1} \cap \Omega_{n2}) \\
&\geq \left(1 - \frac{2}{\pi} \arctan(\zeta_*^{-1}(2a_0 + 24\zeta_0))\right) (1 - 4 \exp(-\log p)) \\
&\geq \left(1 - \frac{2}{\pi} \arctan(\zeta_*^{-1}(2a_0 + 24\zeta_0))\right) / 5.
\end{aligned}$$

The result of the lemma holds with the ζ_1 defined in Lemma D.3 and $\zeta_2 = (1 - \frac{2}{\pi} \arctan(\zeta_*^{-1}(2a_0 + 24\zeta_0))) / 5$. \blacksquare

APPENDIX C. Proofs of the Main Theory

We provide below proofs of Theorems 1–4. Proofs of other results are given in Appendix D.

To prove Theorem 1, we will first establish that the lower bound stated in the theorem holds with high probability for

$$U_n(\Delta) = n^{-1} \sum_{i=1}^n X_i^T \Delta [\mathbf{I}(\epsilon_i \leq X_i^T \Delta) - \mathbf{I}(\epsilon_i \leq 0)], \quad (\text{C.1})$$

which is $\langle S_n(\beta^* + \Delta) - S_n(\beta^*), \Delta \rangle$ corresponding to a specific choice of subgradient in $\partial Q_n(\beta)$. Note that $U_n(\Delta)$ is not Lipschitz continuous in Δ . We start by constructing a Lipschitz continuous lower bound of $U_n(\Delta)$.

We observe that $X_i^T \Delta$ and $\mathbf{I}(\epsilon_i \leq X_i^T \Delta) - \mathbf{I}(\epsilon_i \leq 0)$ always have the same sign. Furthermore, if $X_i^T \Delta > 0$, $\mathbf{I}(\epsilon_i \leq X_i^T \Delta) - \mathbf{I}(\epsilon_i \leq 0)$ is nonzero only $0 < \epsilon_i \leq X_i^T \Delta$; similarly, if

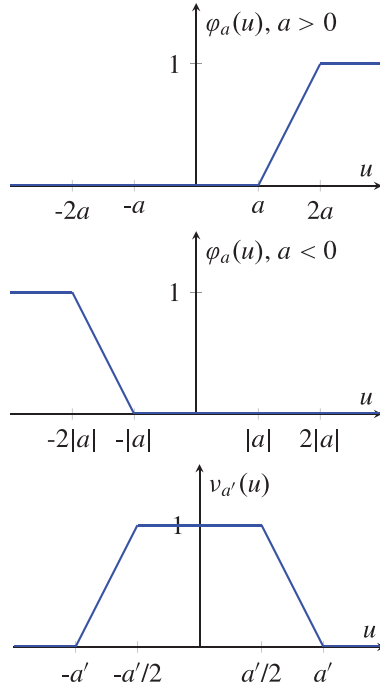
$X_i^T \Delta < 0$, $I(\epsilon_i \leq X_i^T \Delta) - I(\epsilon_i \leq 0)$ is nonzero only if $X_i^T \Delta < \epsilon_i \leq 0$. Hence,

$$\begin{aligned} U_n(\Delta) &= n^{-1} \sum_{i=1}^n X_i^T \Delta I(0 < \epsilon_i \leq X_i^T \Delta) - n^{-1} \sum_{i=1}^n X_i^T \Delta I(X_i^T \Delta < \epsilon_i \leq 0) \\ &\geq n^{-1} \sum_{i=1}^n |\epsilon_i| [I(0 < \epsilon_i < X_i^T \Delta) + I(X_i^T \Delta < \epsilon_i \leq 0)] \\ &\geq n^{-1} \sum_{i=1}^n |\epsilon_i| \varphi_{\epsilon_i}(X_i^T \Delta) v_{b\|\Delta\|_2}(X_i^T \Delta), \end{aligned}$$

where b is a positive constant, and

$$\varphi_a(u) = \begin{cases} 1, & \text{if } |u| > 2|a| \text{ and } ua > 0, \\ -1 + \frac{|u|}{a}, & \text{if } |a| \leq |u| \leq 2|a| \text{ and } ua > 0, \\ 0, & \text{otherwise,} \end{cases}$$

$$v_{a'}(u) = \begin{cases} 1, & \text{if } |u| < \frac{a'}{2}, \\ 2 - \frac{2|u|}{a'}, & \text{if } 0 < \frac{a'}{2} \leq |u| \leq a', \\ 0, & \text{otherwise,} \end{cases}$$



for some $a \in \mathbb{R}$ and $a' \geq 0$. Note that $a\varphi_a(u)$ and $a'v_{a'}(u)$ are both Lipschitz in u : $|a\varphi_a(u_1) - a\varphi_a(u_2)| \leq |u_1 - u_2|$ and $|a'v_{a'}(u_1) - a'v_{a'}(u_2)| \leq 2|u_1 - u_2|$. We also note that

$$I(u > 2a > 0) + I(u < 2a < 0) \leq \varphi_a(u) \leq I(|u| > |a|), \quad (\text{C.2})$$

$$v_{a'}(u) \leq I(|u| \leq a'), \quad (\text{C.3})$$

$$1 - v_{a'}(u) \leq I(|u| \geq a'/2). \quad (\text{C.4})$$

Define

$$V_{nb}(\Delta) = n^{-1} \sum_{i=1}^n |\epsilon_i| \varphi_{\epsilon_i}(X_i^T \Delta) v_{b||\Delta||_2}(X_i^T \Delta).$$

Then $U_n(\Delta) \geq V_{nb}(\Delta)$, $\forall b > 0$. As Lemma C.1 suggests, $V_{nb}(\Delta)$ is a Lipschitz-continuous lower bound of $U_n(\Delta)$.

LEMMA C.1. *Let $g_\epsilon(\xi) = |\epsilon| \varphi_\epsilon(\xi) v_{b\delta}(\xi)$, where $\delta \geq 0$ is a constant. Then $g_\epsilon(\xi)$ is Lipschitz-continuous in ξ ,*

$$|g_\epsilon(\xi) - g_\epsilon(\xi')| \leq 3|\xi - \xi'|. \quad (\text{C.5})$$

Proof. Case I: When $\delta = 0$, we have $v_{b\delta}(\xi) = v_{b\delta}(\xi') = 0$, and (C.5) is satisfied. So we only need to consider the scenario where $\delta > 0$.

Case II: $|\xi| > b\delta > 0$ and $|\xi'| > b\delta > 0$. In this case, $v_{b\delta}(\xi) = v_{b\delta}(\xi') = 0$. Hence, (C.5) holds trivially.

Case III: Assume $\delta > 0$, but case II is not satisfied. Then at least one of ξ or ξ' has absolute value no larger than $b\delta$. Without loss of generality, we assume $|\xi| \leq b\delta$. Then

$$\begin{aligned} |g_\epsilon(\xi) - g_\epsilon(\xi')| &\leq \underbrace{|\epsilon| \varphi_\epsilon(\xi)}_{\leq |\xi| \leq b\delta} \underbrace{|v_{b\delta}(\xi) - v_{b\delta}(\xi')|}_{\leq \frac{2}{b\delta} |\xi - \xi'|} + \underbrace{|\epsilon| |\varphi_\epsilon(\xi) - \varphi_\epsilon(\xi')|}_{\leq |\xi - \xi'|} \underbrace{v_{b\delta}(\xi')}_{\leq 1} \\ &\leq 3|\xi - \xi'|. \end{aligned}$$

Hence, (C.5) holds. On the other hand, if $|\xi'| \leq \tau\delta$, we can bound $|g_\epsilon(\xi) - g_\epsilon(\xi')|$ by $|\epsilon| |\varphi_\epsilon(\xi) - \varphi_\epsilon(\xi')| v_{b\delta}(\xi) + |\epsilon| \varphi_\epsilon(\xi') |v_{b\delta}(\xi) - v_{b\delta}(\xi')|$ and the conclusion follows similarly. \blacksquare

By an application of the dominated convergence theorem, we can prove that there exists a positive constant b_0 such that $\forall b \geq b_0$,

$$\mathbb{E}_{X_i} \left\{ (X_i^T \Delta)^2 \mathbf{1}(|X_i^T \Delta| \geq b||\Delta||_2/2) \right\} \leq \frac{1}{2} \mathbb{E} \{ (X_i^T \Delta)^2 \}. \quad (\text{C.6})$$

LEMMA C.2. *Assume conditions (C1) and (C2) hold, and let b_0 be the aforementioned positive constant. Then there exists a positive constant k_0 such that*

$$E(V_{nb}(\Delta)) \geq a^* ||\Delta||_2^2,$$

uniformly on $\{||\Delta||_2 \leq 1\} \cap \mathbb{C}$, $\forall b \geq b_0$.

Proof. We will first prove (C.6). Write

$$\begin{aligned} & \mathbb{E}_{X_i} \left\{ (X_i^T \Delta)^2 \mathbf{I}(|X_i^T \Delta| \geq b \|\Delta\|_2 / 2) \right\} \\ &= \mathbb{E} \left\{ (X_i^T \Delta)^2 \right\} \mathbb{E}_{X_i} \left\{ \frac{(X_i^T \Delta)^2}{\mathbb{E} \left\{ (X_i^T \Delta)^2 \right\}} \mathbf{I}(|X_i^T \Delta| \geq b \|\Delta\|_2 / 2) \right\}. \end{aligned}$$

We first note that $\mathbb{E} \left\{ \frac{(X_i^T \Delta)^2}{\mathbb{E} \left\{ (X_i^T \Delta)^2 \right\}} \right\} \leq 1$. Furthermore,

$$\begin{aligned} & P \left\{ \frac{(X_i^T \Delta)^2}{\mathbb{E} \left\{ (X_i^T \Delta)^2 \right\}} \mathbf{I}(|X_i^T \Delta| \geq b \|\Delta\|_2 / 2) \neq 0 \right\} \\ &\leq P \left\{ |X_i^T \Delta| \geq b \|\Delta\|_2 / 2 \right\} \\ &\leq \frac{4 \mathbb{E} \left\{ (X_i^T \Delta)^2 \right\}}{b^2 \|\Delta\|_2^2} \leq \frac{4k_u}{b^2}. \end{aligned}$$

By the dominated convergence theorem, $\mathbb{E}_{X_i} \left\{ \frac{(X_i^T \Delta)^2}{\mathbb{E} \left\{ (X_i^T \Delta)^2 \right\}} \mathbf{I}(|X_i^T \Delta| \geq \frac{b \|\Delta\|_2}{2}) \right\} \rightarrow 0$ as $b \rightarrow \infty$. Hence, there exists a positive constant b_0 , whose choice only depends on the probability distributions of X_i , such that $\forall b \geq b_0$, we have

$$\mathbb{E}_{X_i} \left\{ \frac{(X_i^T \Delta)^2}{\mathbb{E} \left\{ (X_i^T \Delta)^2 \right\}} \mathbf{I} \left(|X_i^T \Delta| \geq \frac{b \|\Delta\|_2}{2} \right) \right\} \leq 1/2. \quad (\text{C.7})$$

We have

$$\begin{aligned} & \mathbb{E}(V_{nb}(\Delta)) \\ &= n^{-1} \sum_{i=1}^n \mathbb{E} \left\{ |\epsilon_i| \varphi_{\epsilon_i}() X_i^T \Delta \nu_{b \|\Delta\|_2} (X_i^T \Delta) \right\} \\ &\geq n^{-1} \sum_{i=1}^n \mathbb{E}_{X_i} \left\{ \mathbb{E}_{\epsilon_i | X_i} [|\epsilon_i| \{ \mathbf{I}(X_i^T \Delta > 2\epsilon_i > 0) + \mathbf{I}(X_i^T \Delta < 2\epsilon_i < 0) \}] \mathbf{I}(|X_i^T \Delta| \leq b \|\Delta\|_2 / 2) \right\} \\ &= n^{-1} \sum_{i=1}^n \mathbb{E}_{X_i} \left\{ \left(\int_0^{|X_i^T \Delta|/2} t f_i(t) dt + \int_{-|X_i^T \Delta|/2}^0 (-t) f_i(t) dt \right) \mathbf{I}(|X_i^T \Delta| \leq b \|\Delta\|_2 / 2) \right\} \\ &= n^{-1} \sum_{i=1}^n \mathbb{E}_{X_i} \left\{ \left(\int_0^{|X_i^T \Delta|/2} t [f_i(t) + f_i(-t)] dt \right) \mathbf{I}(|X_i^T \Delta| \leq b \|\Delta\|_2 / 2) \right\} \\ &\geq \frac{m_0}{4} \mathbb{E}_{X_i} \left\{ (X_i^T \Delta)^2 \mathbf{I}(|X_i^T \Delta| \leq b \|\Delta\|_2 / 2) \right\} \\ &= \frac{m_0}{4} \mathbb{E}_{X_i} \left\{ (X_i^T \Delta)^2 \right\} - \frac{m_0}{4} \mathbb{E}_{X_i} \left\{ (X_i^T \Delta)^2 \mathbf{I}(|X_i^T \Delta| > b \|\Delta\|_2 / 2) \right\} \\ &\geq \frac{m_0}{8} \mathbb{E}_{X_i} \left\{ (X_i^T \Delta)^2 \right\} \\ &\geq a^* \|\Delta\|_2^2, \end{aligned}$$

where $a^* = (m_0 m_1)/8$, the first inequality applies (C.2) and (C.4), the second inequality applies condition (C1), the second last inequality applies (C.6), and the last inequality follows from condition (C2). \blacksquare

For an arbitrary $0 < \delta \leq 1$, let $S_2(\delta) = \{\Delta : \|\Delta\|_2 = \delta\}$. Let $\Gamma(t) = \{\Delta : \Delta \in S_2(\delta), \|\Delta\|_1 \leq t\|\Delta\|_2\} \cap \mathbb{C}$, for an arbitrary $t > 0$. Lemma C.3 establishes a concentration inequality for

$$V_{nb}(\Delta) = n^{-1} \sum_{i=1}^n |\epsilon_i| \varphi_{\epsilon_i}(X_i^T \Delta) v_b \|\Delta\|_2 (X_i^T \Delta),$$

where $b \geq b_0$ is a positive constant. Consider the event $A_{n1} = \{\max_{1 \leq j \leq p} \hat{\sigma}_j^2 \leq m_x\}$, where $\hat{\sigma}_j^2 = n^{-1} \sum_{i=1}^n x_{ij}^2$. Then $P(A_{n1}) \geq 1 - \delta_n$, where $\delta_n \rightarrow 0$ as $n \rightarrow \infty$.

LEMMA C.3. *For an arbitrary $b \geq b_0$, define*

$$Z_n(t) = \sup_{\Delta \in \Gamma(t)} |V_{nb}(\Delta) - E(V_{nb}(\Delta))|. \quad (\text{C.8})$$

Then there exists a positive constant c^ which does not depend on (n, p, t) such that on the event A_{n1} ,*

$$P(Z_n(t) \geq c^* t \delta \sqrt{\log p / n}) \leq \exp\left(-\frac{c^{*2} t^2}{16b^2} \log p\right).$$

Proof. First, we note that by (C.2) and (C.3), $\forall \Delta \in \Gamma(t)$,

$$0 \leq |\epsilon_i| \varphi_{\epsilon_i}(X_i^T \Delta) v_b \|\Delta\|_2 (X_i^T \Delta) \leq |\epsilon_i| \mathbf{I}(|X_i^T \Delta| > |\epsilon_i|) \mathbf{I}(|X_i^T \Delta| \leq b\delta) \leq b\delta.$$

If we change one observation of the sample, the value of $Z_n(t)$ changes at most $2b\delta/n$. By the bounded difference inequality, $\forall s > 0$,

$$P(Z_n(t) - E(Z_n(t)) \geq s) \leq \exp\left(-\frac{ns^2}{4b^2\delta^2}\right). \quad (\text{C.9})$$

Next, we evaluate $E(Z_n(t))$. Let $\{\pi_1, \dots, \pi_n\}$ be a Rademacher sequence of Bernoulli random variables independent of (X_i, ϵ_i) . Let $g_\epsilon(\xi) = |\epsilon| \varphi_\epsilon(\xi) v_b \delta(\xi)$ be the function defined in Lemma C.1. We have

$$\begin{aligned} E(Z_n(t)) &\leq 2E_{(\pi, X, \epsilon)} \left(\sup_{\Delta \in \Gamma(t)} n^{-1} \sum_{i=1}^n \pi_i g_{\epsilon_i}(X_i^T \Delta) \right) \quad (\text{symmetrization}) \\ &\leq 12E_{(X, \epsilon)} E_{\pi|X, \epsilon} \left(\sup_{\Delta \in \Gamma(t)} n^{-1} \sum_{i=1}^n \pi_i X_i^T \Delta \right) \quad (\text{contraction}) \\ &\leq 12 \sup_{\Delta \in \Gamma(t)} \|\Delta\|_1 E_X E_{\pi|X} \left(n^{-1} \left\| \sum_{i=1}^n \pi_i X_i \right\|_\infty \right) \\ &\leq 12t\delta E_X E_{\pi|X} \left(n^{-1} \left\| \sum_{i=1}^n \pi_i X_i \right\|_\infty \right), \end{aligned} \quad (\text{C.10})$$

where the first inequality follows from the symmetrization theorem (van der Vaart and Wellner, 1996), and the second inequality applies Lemma C.1 and the contraction theorem (Ledoux and Talagrand, 2013). Next, we will evaluate $E_{\pi|X}(n^{-1} \|\sum_{i=1}^n \pi_i X_i\|_\infty)$. We

observe that, for $1 \leq j \leq p$, conditional on X_i , $n^{-1} \sum_{i=1}^n \pi_i x_{ij}$ is mean-zero sub-Gaussian with parameter bounded by $n^{-1} (\sum_{i=1}^n x_{ij}^2)^{1/2}$. By the property of sub-Gaussian random variables (van de Geer, 2016, Lem. 17.5), we have

$$\begin{aligned} E_{\pi|X} \left(\max_{1 \leq j \leq p} n^{-1} \left| \sum_{i=1}^n \pi_i x_{ij} \right| \right) &\leq \sqrt{2 \log(2p)/n} \max_{1 \leq j \leq p} \left(n^{-1} \sum_{i=1}^n x_{ij}^2 \right)^{1/2} \\ &\leq \sqrt{2 m_x \log(2p)/n} \end{aligned}$$

on the event A_{n1} . Hence, on A_{n1} ,

$$E(Z_n(t)) \leq \frac{1}{2} c^* t \delta \sqrt{\log p/n}, \quad (\text{C.11})$$

where $c^* = 48 \sqrt{m_x}$. Taking $s = \frac{1}{2} c^* t \delta \sqrt{\log p/n}$ in (C.9), we have that on A_{n1} ,

$$P(Z_n(t) \geq c^* t \delta \sqrt{\log p/n}) \leq \exp \left(- \frac{c^{*2} t^2}{16 b^2} \log p \right). \quad \blacksquare$$

Proof of Theorem 1. The proof consists of two steps. At the first step, we will establish that the lower bound stated in the theorem holds with high probability for $U_n(\Delta)$ defined in (C.1), which is $\langle S_n(\beta^* + \Delta) - S_n(\beta^*), \Delta \rangle$ corresponding to a specific choice of subgradient in $\partial Q_n(\beta)$. At the second step, we will show that the lower bound holds for an arbitrary choice of subgradient in $\partial Q_n(\beta)$.

Step 1. Note that $U_n(\Delta) \geq V_{nb}(\Delta)$, $\forall b > 0$. We will first prove that, for an arbitrary $b \geq b_0$, there exist some positive constants a_1 and a_2 such that with probability at least $1 - \delta_n - a_1 \exp(-a_2 \log p)$,

$$V_{nb}(\Delta) \geq a^* \|\Delta\|_2^2 - c^* \|\Delta\|_1 \sqrt{\frac{\log p}{n}} \quad (\text{C.12})$$

uniformly over all $\Delta \in S_2(\delta) \cap \mathbb{C}$, where c^* is the positive constant in Lemma C.3, and a_1 and a_2 (specified later in this proof) only depend on the probability distributions of X_i and ϵ_i . Recall that $S_2(\delta) = \{\Delta : \|\Delta\|_2 = \delta\}$, $0 < \delta \leq 1$. Since $0 < \delta \leq 1$, we have $\frac{\|\Delta\|_1}{\|\Delta\|_2} \geq 1$. Noting that $V_n(\Delta)$ is always nonnegative, we only need to verify (C.12) for Δ satisfying $a^* \|\Delta\|_2^2 - c^* \|\Delta\|_1 \sqrt{\frac{\log p}{n}} \geq 0$. It is sufficient to consider $\Delta \in \tilde{S}_2(\delta)$, where

$$\tilde{S}_2(\delta) = \left\{ \Delta : \Delta \in S_2(\delta) \cap \mathbb{C}, 1 \leq \frac{\|\Delta\|_1}{\|\Delta\|_2} \leq \frac{a^*}{c^*} \delta \sqrt{\frac{n}{\log p}} \right\}.$$

Lemma C.2 ensures that $E(V_{nb}(\Delta)) \geq a^* \|\Delta\|_2^2$ for any $\Delta \in \tilde{S}_2(\delta)$. We verify (C.12) by proving

$$|V_{nb}(\Delta) - E(V_{nb}(\Delta))| \leq c^* \delta \frac{\|\Delta\|_1}{\|\Delta\|_2} \sqrt{\log p/n} \quad (\text{C.13})$$

holds uniformly for $\Delta \in \tilde{S}_2(\delta)$ with probability at least $1 - a_1 \exp(-a_2 \log p)$ for some positive constants a_1 and a_2 . We will prove (C.13) by employing a peeling technique (see, e.g., van der Vaart and Wellner, 1996; van de Geer, 2000, and the references therein). Write

$$h(\Delta) = \frac{\|\Delta\|_1}{\|\Delta\|_2} \quad \text{and} \quad g(h(\Delta)) = c^* \delta h(\Delta) \sqrt{\log p/n}.$$

Define

$$B_m = \{\Delta : 2^{m-1}\mu \leq g(h(\Delta)) \leq 2^m\mu\} \cap \tilde{S}_2(\delta), \quad m = 1, \dots, M,$$

where $\mu = c^*\delta\sqrt{\log p/n}$, and M is taken as the smallest positive integer such that $2^M \geq \frac{a^*\delta}{c^*}\sqrt{n/\log p}$. Note that if $\Delta \in B_m$, then $2^{m-1} \leq h(\Delta) \leq 2^m$, $m = 1, \dots, M$.

Let $E = \{\text{the event in (C.13) holds for all } \Delta \in \tilde{S}_2(\delta)\}$. Recall the event $A_{n1} = \{\max_{1 \leq j \leq p} \hat{\sigma}_j^2 \leq m_x\}$, where $\hat{\sigma}_j^2 = n^{-1} \sum_{i=1}^n x_{ij}^2$. Then, on the event A_{n1} , we have

$$\begin{aligned} P(E^c) &\leq \sum_{m=1}^M P\left(\exists \Delta \in \tilde{S}_2(\delta) \cap B_m \text{ such that } |V_{nb}(\Delta) - E(V_{nb}(\Delta))| > g(h(\Delta))\right) \\ &\leq \sum_{m=1}^M P\left(\sup_{\Delta \in \tilde{S}_2(\delta) \cap \mathbb{C}, 2^{m-1} \leq h(\Delta) \leq 2^m} |V_{nb}(\Delta) - E(V_{nb}(\Delta))| > c^*\delta 2^{m-1} \sqrt{\log p/n}\right) \\ &\leq \sum_{m=1}^M \exp\left(-a_2 2^{2m-2} \log p\right) \quad (\text{by Lemma C.3, taking } a_2 = c^{*2}/(16b^2)) \\ &\leq \sum_{m=1}^M \exp\left(-ma_2 \log p\right) \quad (\text{noting } m \leq 2^{2m-2}, m = 1, 2, \dots) \\ &\leq \exp(-a_2 \log p) \sum_{m=0}^{\infty} (\exp(-a_2 \log p))^m \\ &= \frac{\exp(-a_2 \log p)}{1 - \exp(-a_2 \log p)} \quad (\text{sum of a geometric series}) \\ &\leq a_1 \exp(-a_2 \log p), \end{aligned}$$

where $a_1 = 1 + \exp(-a_2 \log 2)$, and the third inequality applies Lemma C.3 on the event A_{n1} . Note that $P(A_{n1}) \geq 1 - \delta_n$, for $\delta_n \rightarrow 0$. Hence, with probability at least $1 - \delta_n - a_1 \exp(-a_2 \log p)$, $U_n(\Delta) \geq a^* \|\Delta\|_2^2 - c^* \|\Delta\|_1 \sqrt{\frac{\log p}{n}}$ uniformly over all $\Delta \in \tilde{S}_2(\delta) \cap \mathbb{C}$.

Step 2. We now consider an arbitrary subgradient $S_n(\beta) = (S_{n1}(\beta), \dots, S_{np}(\beta))^T$ in $\partial Q_n(\beta)$.

According to (6), $S_j(\beta)$, $j = 1, \dots, p$, has the form

$$\begin{aligned} S_{nj}(\beta) &= -\tau n^{-1} \sum_{i=1}^n x_{ij} \beta + n^{-1} \sum_{i=1}^n x_{ij} \mathbf{I}(Y_i - X_i^T \beta \leq 0) \\ &\quad - n^{-1} \sum_{i=1}^n [v_i + (1 - \tau)] x_{ij} \mathbf{I}(Y_i - X_i^T \beta = 0), \end{aligned}$$

where $v_i = 0$ if $Y_i - X_i^T \beta \neq 0$, and $v_i \in [\tau - 1, \tau]$; otherwise, $i = 1, \dots, n$. Hence,

$$\begin{aligned} &\langle S_n(\beta^* + \Delta) - S_n(\beta^*), \Delta \rangle \\ &\geq n^{-1} \sum_{i=1}^n X_i^T \Delta [\mathbf{I}(\epsilon_i \leq X_i^T \Delta) - \mathbf{I}(\epsilon_i \leq 0)] - n^{-1} \sum_{i=1}^n |X_i^T \Delta| \mathbf{I}(\epsilon_i = X_i^T \Delta) \end{aligned}$$

$$\begin{aligned}
& -n^{-1} \sum_{i=1}^n |X_i^T \Delta| \mathbf{I}(\epsilon_i = 0) \\
& = U_n(\Delta) - n^{-1} \sum_{i=1}^n |X_i^T \Delta| \mathbf{I}(\epsilon_i = X_i^T \Delta) - n^{-1} \sum_{i=1}^n |X_i^T \Delta| \mathbf{I}(\epsilon_i = 0).
\end{aligned}$$

Consider any Δ such that $\|\Delta\|_2 \leq 1$. Consider the event $G_n = \{U_n(\Delta) < k_0 \|\Delta\|_2^2 - c^* \|\Delta\|_1 \sqrt{\log p/n}\}$. We have

$$\begin{aligned}
& P\left(\langle S_n(\beta) - S_n(\beta^*), \Delta \rangle < k_0 \|\Delta\|_2^2 - c^* \|\Delta\|_1 \sqrt{\log p/n}\right) \\
& \leq P\left(U_n(\Delta) - n^{-1} \sum_{i=1}^n |X_i^T \Delta| (\mathbf{I}(\epsilon_i = X_i^T \Delta) + \mathbf{I}(\epsilon_i = 0)) \right. \\
& \quad \left. < k_0 \|\Delta\|_2^2 - c^* \|\Delta\|_1 \sqrt{\log p/n}, G_n^c\right) + P(G_n) \\
& \leq P\left(n^{-1} \sum_{i=1}^n |X_i^T \Delta| (\mathbf{I}(\epsilon_i = X_i^T \Delta) + \mathbf{I}(\epsilon_i = 0)) > 0\right) + P(G_n) \\
& \leq 0 + \delta_n + a_1 \exp(-a_2 \log p) = \delta_n + a_1 \exp(-a_2 \log p),
\end{aligned}$$

where the last inequality follows because ϵ_i is assumed to have a continuous density (e.g., Ruppert and Carroll, 1980, Lem. A.1). This proves the theorem. \blacksquare

We first present a lemma that shows that the uniform lower bound in Theorem 1 on $\{\|\Delta\|_2 \leq 1\} \cap \mathbb{C}$ can be extended to $\{\|\Delta\|_2 > 1\} \cap \mathbb{C}$.

LEMMA C.4. *Suppose conditions (C1)–(C3) are satisfied. There exist some positive constants a^* , c^* , a_1 , and a_2 , such that for any subgradient $S_n \in \partial Q_n(\beta)$, with probability at least $1 - \delta_n - a_1 \exp(-a_2 \log p)$,*

$$\langle S_n(\beta) - S_n(\beta^*), \Delta \rangle \geq a^* \|\Delta\|_2 - c^* \sqrt{\frac{\log p}{n}} \|\Delta\|_1, \quad (\text{C.14})$$

uniformly on $\{\|\Delta\|_2 > 1\} \cap \mathbb{C}$, where the positive constants a^ , c^* , a_1 , and a_2 are those defined in Theorem 1.*

Proof. We first observe that, $\forall 0 < s < 1, \forall t \in \mathbb{R}$,

$$t[\mathbf{I}(\epsilon_i \leq t) - \mathbf{I}(\epsilon_i \leq 0)] \geq t[\mathbf{I}(\epsilon_i \leq st) - \mathbf{I}(\epsilon_i \leq 0)] \geq 0.$$

We thus have, $\forall 0 < s < 1$,

$$\begin{aligned}
U_n(\Delta) &= n^{-1} \sum_{i=1}^n X_i^T \Delta [\mathbf{I}(\epsilon_i \leq X_i^T \Delta) - \mathbf{I}(\epsilon_i \leq 0)] \\
&\geq \frac{1}{s} n^{-1} \sum_{i=1}^n X_i^T (s\Delta) [\mathbf{I}(\epsilon_i \leq X_i^T s\Delta) - \mathbf{I}(\epsilon_i \leq 0)] \\
&= \frac{1}{s} U_n(s\Delta).
\end{aligned} \quad (\text{C.15})$$

For $\|\Delta\|_2 > 1$, we take $s = \frac{1}{\|\Delta\|_2} \in (0, 1)$ in (C.15) to rescale Δ to $\frac{\Delta}{\|\Delta\|_2}$ and obtain

$$\begin{aligned} U_n(\Delta) &\geq \|\Delta\|_2 U_n\left(\frac{\Delta}{\|\Delta\|_2}\right) \\ &\geq \|\Delta\|_2 \left(a^* - c^* \frac{\|\Delta\|_1}{\|\Delta\|_2} \sqrt{\frac{\log p}{n}}\right) \\ &\geq a^* \|\Delta\|_2 - c^* \|\Delta\|_1 \sqrt{\frac{\log p}{n}}, \end{aligned} \quad (\text{C.16})$$

with probability at least $1 - \delta_n - a_1 \exp(-a_2 \log p)$, where the second inequality applies Theorem 1. Furthermore, arguing the same way as in the proof of Theorem 1, we can show that the above lower bound holds for an arbitrary subgradient $S_n(\beta)$ in $\partial Q_n(\beta)$. ■

Proof of Theorem 2 (QR-LASSO). Let $\hat{\beta}$ be the QR-LASSO estimator. Then there exists some subgradient $S_n(\beta) \in \partial Q_n(\beta)$ such that

$$(S_n(\hat{\beta}) + \lambda \text{sgn}(\hat{\beta}), \beta^* - \hat{\beta}) = 0, \quad (\text{C.17})$$

where $\text{sgn}(\hat{\beta}) = (0, \text{sgn}(\hat{\beta}_2), \dots, \text{sgn}(\hat{\beta}_p))^T$.

Let $\hat{v} = \hat{\beta} - \beta^*$. We will first use proof by contradiction to show that

$$P(\|\hat{v}\|_2 \leq 1) \geq 1 - 2\delta_n - 2\exp(-\log p) - a_1 \exp(-a_2 \log p).$$

Let $\Lambda_n = \{\lambda \geq 2\|\tilde{S}_n\|_\infty\}$. The choice of λ and Lemma 3 imply that $P(\Lambda_n) \geq 1 - \delta_n - 2\exp(-\log p)$. Assume $\|\hat{v}\|_2 > 1$. Define the event $B_{n1} = \{(S_n(\hat{\beta}) - S_n(\beta^*), \hat{v}) \geq a^* \|\hat{v}\|_2 - c^* \|\hat{v}\|_1 \sqrt{\log p/n}\}$. If $\|\hat{v}\|_2 > 1$, Lemma C.4 implies that $P(B_{n1}) \geq 1 - \delta_n - a_1 \exp(-a_2 \log p)$. It follows that $P(\Lambda_n \cap B_{n1}) \geq 1 - 2\delta_n - 2\exp(-\log p) - a_1 \exp(-a_2 \log p)$. It is sufficient to show that a contradiction occurs on the event $\Lambda_n \cap B_{n1}$.

To see the contradiction, we first observe that on B_{n1} , (C.17) implies that

$$(-\lambda \text{sgn}(\hat{\beta}) - S_n(\beta^*), \hat{v}) \geq a^* \|\hat{v}\|_2 - c^* \|\hat{v}\|_1 \sqrt{\log p/n}. \quad (\text{C.18})$$

On the other hand, by Hölder's inequality,

$$(-\lambda \text{sgn}(\hat{\beta}) - S_n(\beta^*), \hat{v}) \leq \{\lambda + \|S_n(\beta^*)\|_\infty\} \|\hat{v}\|_1.$$

As the argument in Step 2 of the proof of Theorem 1 implies, for an arbitrary subgradient $S_n(\cdot)$, on the event Λ_n , $\|S_n(\beta^*)\|_\infty \leq 2\lambda$ with probability one. Hence, on Λ_n ,

$$(-\lambda \text{sgn}(\hat{\beta}) - S_n(\beta^*), \hat{v}) \leq 1.5k_0 \sqrt{\log p/n} \|\hat{v}\|_1. \quad (\text{C.19})$$

(C.18) and (C.19) together imply that

$$a^* \|\hat{v}\|_2 \leq (2k_0 + c^*) \sqrt{\log p/n} \|\hat{v}\|_1. \quad (\text{C.20})$$

(i) (Hard-sparsity case) The proof of Lemma 2 ensures that on the event Λ_n , we have $\hat{v} \in \Gamma_H$ and hence $\|\hat{v}\|_1 \leq 4\sqrt{s} \|\hat{v}\|_2$. Then, by (C.20), $a^* \leq 4(2k_0 + c^*) \sqrt{s \log p/n}$. This contradicts with the assumption $n > (a_1^*)^2 s \log p$, where $a_1^* = 4(2k_0 + c^*)/a^*$.

(ii) (Soft sparsity case) The proof of Lemma 2 ensures that on the event Λ_n , we have $\hat{v} \in \Gamma_W = \{v \in \mathbb{R}^p : \|v_{S_a^c}\|_1 \leq 3\|v_{S_a}\|_1 + 4\|\beta_{S_a^c}^*\|_1\}$, where $S_a = S_{-a} \cup \{1\}$ with the index set $S_{-a} = \{j : |\beta_j^*| > a, 2 \leq j \leq p\}$ and S_a^c denotes the complement of S_a in $\{1, 2, \dots, p\}$. Under the soft sparsity assumption, we have $a|S_{-a}| \leq \sum_{i=2}^p |\beta_i^*| \leq R$, where $|S_{-a}|$ denotes

the cardinality of the set S_{-a} . Hence, $|S_{-a}| \leq a^{-1}R$. Since $\widehat{v} \in \Gamma_W$, we have

$$\|\widehat{v}\|_1 \leq 4(\|\widehat{v}_{S_a}\|_1 + \|\beta_{S_a}^*\|_1) \leq 4(\sqrt{a^{-1}R+1}\|\widehat{v}\|_2 + R), \quad (\text{C.21})$$

which holds for any $a > 0$. Taking $a = R/3$, we obtain

$$\|\widehat{v}\|_1 \leq 4[2\|\widehat{v}\|_2 + R]. \quad (\text{C.22})$$

Hence, (C.20) contradicts with the assumption $a_1^* \sqrt{\log p/n} \max\{2, R\} < 1/2$.

Define the event $D_{n1} = \{\|\widehat{v}\|_2 \leq 1\}$. In the above, we have verified that $P(D_{n1}) \geq 1 - 2\delta_n - 2\exp(-\log p) - a_1 \exp(-a_2 \log p)$. Define $B_{n2} = \{S_n(\widehat{\beta}) - S_n(\beta^*), \widehat{\beta} - \beta^*\} \geq a^* \|\widehat{v}\|_2^2 - c^* \|\widehat{v}\|_1 \sqrt{\log p/n}\}$. By Theorem 1, on D_{n1} , $P(B_{n2}) \geq 1 - \delta_n - a_1 \exp(-a_2 \log p)$. From now on, we consider the event $\Lambda_n \cap D_{n1} \cap B_{n1}$. We have $P(\Lambda_n \cap D_{n1} \cap B_{n1}) \geq 1 - 4\delta_n - 4\exp(-\log p) - 2a_1 \exp(-a_2 \log p)$. By the convexity of $\|\beta\|_1$, we have

$$\|\beta^*\|_1 - \|\widehat{\beta}\|_1 \geq \langle \text{sgn}(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle. \quad (\text{C.23})$$

Combining (C.23) with (C.17), we have

$$\langle S_n(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle = -\langle \lambda \text{sgn}(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle \geq \lambda(\|\widehat{\beta}\|_1 - \|\beta^*\|_1). \quad (\text{C.24})$$

On B_{n2} , (C.24) implies that

$$\langle -S_n(\beta^*), \widehat{\beta} - \beta^* \rangle \geq a^* \|\widehat{v}\|_2^2 + \lambda(\|\widehat{\beta}\|_1 - \|\beta^*\|_1) - c^* \|\widehat{v}\|_1 \sqrt{\log p/n}.$$

Applying Hölder's inequality, we have

$$a^* \|\widehat{v}\|_2^2 + \lambda(\|\widehat{\beta}\|_1 - \|\beta^*\|_1) - c^* \|\widehat{v}\|_1 \sqrt{\log p/n} \leq \|S_n(\beta^*)\|_\infty \|\widehat{v}\|_1.$$

Rearranging terms and applying the triangle inequality, we obtain

$$\begin{aligned} a^* \|\widehat{v}\|_2^2 &\leq (\lambda + c^* \sqrt{\log p/n} + \|S_n(\beta^*)\|_\infty) \|\widehat{v}\|_1 \\ &\leq (2k_0 + c^*) \sqrt{\log p/n} \|\widehat{v}\|_1. \end{aligned} \quad (\text{C.25})$$

(i) For the hard sparsity case, on Λ_n , $\|\widehat{v}\|_1 \leq 4\sqrt{s}\|\widehat{v}\|_2$. Then (C.25) implies

$$\|\widehat{v}\|_2 \leq a_1^* \sqrt{s \log p/n}.$$

We also obtain $\|\widehat{v}\|_1 \leq 4a_1^* s \sqrt{\log p/n}$.

(ii) For the soft sparsity case, on Λ_n , (C.25) and (C.21) imply that

$$\|\widehat{v}\|_2 \leq 2 \max \left\{ a_1^* \sqrt{a^{-1}R+1} \sqrt{\log p/n}, \sqrt{a_1^* R^{1/2} (\log p/n)^{1/4}} \right\},$$

for any $a > 0$. Taking $a = \sqrt{\log p/n}$, we have

$$\sqrt{a^{-1}R+1} \sqrt{\log p/n} = \sqrt{R + (\log p/n)^{1/2} (\log p/n)^{1/4}} \leq \sqrt{2R} (\log p/n)^{1/4}.$$

We obtain

$$\|\widehat{v}\|_2 \leq a_2^* R^{1/2} (\log p/n)^{1/4},$$

where $a_2^* = 2 \max \left\{ \sqrt{2}a_1^*, \sqrt{a_1^*} \right\}$. The same reasoning that leads to (C.22) also implies that

$$\begin{aligned} \|\widehat{v}\|_1 &\leq 4(\sqrt{|S_a|}\|\widehat{v}\|_2 + \|\beta_{S_a^c}^*\|_1) \\ &\leq 4(a_2^*\sqrt{|S_a|}R^{1/2}(\log p/n)^{1/4} + \|\beta_{S_a^c}^*\|_1), \end{aligned}$$

for any $a > 0$, where $|S_a|$ denotes the cardinality of the set S_a . ■

Proof of Theorem 3 (QR-NCP). The proof is based on the same idea as in the proof of Theorem 2 but is more involved. We consider any feasible local solution $\widehat{\beta}$ such that $\|\widehat{\beta}\|_1 \leq \kappa$ and (8) is satisfied. Then there exists some subgradient $S_n(\beta) \in \partial Q_n(\beta)$ such that

$$\langle S_n(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}) - H'(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle = 0, \quad (\text{C.26})$$

where $H'(\widehat{\beta}) = (0, h'_\lambda(\widehat{\beta}_2), \dots, h'_\lambda(\widehat{\beta}_p))^T$ and $\text{sgn}(\widehat{\beta}) = (0, \text{sgn}(\widehat{\beta}_2), \dots, \text{sgn}(\widehat{\beta}_p))^T$. Denote $\widehat{v} = \widehat{\beta} - \beta^*$.

Let $\widetilde{\Lambda}_n = \{\lambda \geq 4\|\widetilde{S}_n\|_\infty\}$, where λ satisfies the condition of Theorem 3. Lemma 3 implies that $P(\widetilde{\Lambda}_n) \geq 1 - \delta_n - 2\exp(-\log p)$. Let $B_{n1} = \{\langle S_n(\widehat{\beta}) - S_n(\beta^*), \widehat{v} \rangle \geq a^*\|\widehat{v}\|_2 - c^*\|\widehat{v}\|_1\sqrt{\log p/n}\}$. If $\|\widehat{v}\|_2 > 1$, Lemma C.4 implies that $P(B_{n1}) \geq 1 - \delta_n - a_1 \exp(-a_2 \log p)$. We have $P(\widetilde{\Lambda}_n \cap B_{n1}) \geq 1 - 2\delta_n - 2\exp(-\log p) - a_1 \exp(-a_2 \log p)$. We will first use proof by contradiction to show that

$$P(\|\widehat{v}\|_2 \leq 1) \geq 1 - 2\delta_n - 2\exp(-\log p) - a_1 \exp(-a_2 \log p).$$

Assume $\|\widehat{v}\|_2 > 1$. It is sufficient to show that a contradiction occurs on the event $\widetilde{\Lambda}_n \cap B_{n1}$.

On the event B_{n1} , by (C.26), we have

$$\langle H'(\widehat{\beta}) - \lambda \text{sgn}(\widehat{\beta}) - S_n(\beta^*), \widehat{v} \rangle \geq a^*\|\widehat{v}\|_2 - c^*\|\widehat{v}\|_1\sqrt{\frac{\log p}{n}}. \quad (\text{C.27})$$

On the other hand, by Hölder's inequality,

$$\langle H'(\widehat{\beta}) - \lambda \text{sgn}(\widehat{\beta}) - S_n(\beta^*), \widehat{v} \rangle \leq \{\|H'(\widehat{\beta}) - \lambda \text{sgn}(\widehat{\beta})\|_\infty + \|S_n(\beta^*)\|_\infty\}\|\widehat{v}\|_1.$$

Note that $-H(\beta_-) + \lambda\|\beta_- \|_1 = p_\lambda(\beta_-) = \sum_{j=2}^p p_\lambda(|\beta_j|)$. We have $\|H'(\widehat{\beta}) - \lambda \text{sgn}(\widehat{\beta})\|_\infty = \|\partial p_\lambda(\widehat{\beta}_-)\|_\infty$, which is upper bounded by λ (see, e.g., Loh and Wainwright, 2015, Lem. 4). Furthermore, on $\widetilde{\Lambda}_n$, $\|S_n(\beta^*)\|_\infty \leq \lambda/4$. Hence,

$$\langle H'(\widehat{\beta}) - \lambda \text{sgn}(\widehat{\beta}) - S_n(\beta^*), \widehat{v} \rangle \leq \frac{5\lambda}{4}\|\widehat{v}\|_1. \quad (\text{C.28})$$

By the choice of λ , (C.27) and (C.28) together imply that

$$a^*\|\widehat{v}\|_2 \leq (5\lambda/4 + c^*\sqrt{\log p/n})\|\widehat{v}\|_1 \leq \frac{3}{2}\lambda\|\widehat{v}\|_1 \leq \frac{3\kappa k_0}{2}\sqrt{\log p/n}.$$

Hence, we have a contradiction under the assumption $\sqrt{\log p/n} < \frac{2a^*}{3\kappa k_0}$.

Define the event $D_{n1} = \{\|\widehat{v}\|_2 \leq 1\}$. In the above, we have verified that $P(D_{n1}) \geq 1 - 2\delta_n - 2\exp(-\log p) - a_1 \exp(-a_2 \log p)$. Define $B_{n2} = \{\langle S_n(\widehat{\beta}) - S_n(\beta^*), \widehat{\beta} - \beta^* \rangle \geq a^*\|\widehat{v}\|_2^2 - c^*\|\widehat{v}\|_1\sqrt{\log p/n}\}$. By Theorem 1, on D_{n1} , $P(B_{n2}) \geq 1 - \delta_n - a_1 \exp(-a_2 \log p)$. From now on, we consider the event $\widetilde{\Lambda}_n \cap D_{n1} \cap B_{n1}$. We have $P(\widetilde{\Lambda}_n \cap D_{n1} \cap B_{n1}) \geq 1 - 4\delta_n - 4\exp(-\log p) - 2a_1 \exp(-a_2 \log p)$.

By the assumption on the penalty function, $\frac{\gamma_0}{2} \|\beta\|_2^2 - H(\beta_-) + \lambda \|\beta_- \|_1$ is convex in β . This convexity property implies that

$$\begin{aligned} & \left[\frac{\gamma_0}{2} \|\beta^*\|_2^2 - H(\beta_-^*) + \lambda \|\beta_-^*\|_1 \right] - \left[\frac{\gamma_0}{2} \|\widehat{\beta}\|_2^2 - H(\widehat{\beta}_-) + \lambda \|\widehat{\beta}_-\|_1 \right] \\ & \geq \langle \gamma_0 \widehat{\beta} - H'(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle. \end{aligned}$$

So we have

$$\begin{aligned} & p_\lambda(\beta_-^*) - p_\lambda(\widehat{\beta}_-) + \frac{\gamma_0}{2} (\|\beta^*\|_2^2 - \|\widehat{\beta}\|_2^2) \\ & \geq \langle -H'(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle + \gamma_0 \langle \widehat{\beta}, \beta^* \rangle - \gamma_0 \|\widehat{\beta}\|_2. \end{aligned}$$

Or equivalently,

$$p_\lambda(\beta_-^*) - p_\lambda(\widehat{\beta}_-) + \frac{\gamma_0}{2} \|\beta^* - \widehat{\beta}\|_2^2 \geq \langle -H'(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle. \quad (\text{C.29})$$

It follows from (C.26) and (C.29) that

$$\begin{aligned} & \langle S_n(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle \\ & = -\langle -H'(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle \geq p_\lambda(\widehat{\beta}_-) - p_\lambda(\beta_-^*) - \frac{\gamma_0}{2} \|\widehat{\nu}\|_2^2. \end{aligned} \quad (\text{C.30})$$

By Theorem 1, on B_{n2} ,

$$\langle S_n(\widehat{\beta}) - S_n(\beta^*), \widehat{\beta} - \beta^* \rangle \geq a^* \|\widehat{\nu}\|_2^2 - c^* \|\widehat{\nu}\|_1 \sqrt{\log p/n}. \quad (\text{C.31})$$

(C.30) and (C.31) together imply that

$$\langle -S_n(\beta^*), \widehat{\nu} \rangle \geq (a^* - \gamma_0/2) \|\widehat{\nu}\|_2^2 + (p_\lambda(\widehat{\beta}_-) - p_\lambda(\beta_-^*)) - c^* \|\widehat{\nu}\|_1 \sqrt{\log p/n}.$$

Applying Hölder's inequality, we have

$$(a^* - \gamma_0/2) \|\widehat{\nu}\|_2^2 + (p_\lambda(\widehat{\beta}_-) - p_\lambda(\beta_-^*)) - c^* \|\widehat{\nu}\|_1 \sqrt{\log p/n} \leq \|S_n(\beta^*)\|_\infty \|\widehat{\nu}\|_1.$$

Hence,

$$(a^* - \gamma_0/2) \|\widehat{\nu}\|_2^2 \leq (p_\lambda(\beta_-^*) - p_\lambda(\widehat{\beta}_-)) + (c^* \sqrt{\log p/n} + \|S_n(\beta^*)\|_\infty) \|\widehat{\nu}\|_1.$$

On $\widetilde{\Lambda}_n$, we have $\|S_n(\beta^*)\|_\infty \leq \frac{\lambda}{4}$, and the choice of λ implies $c^* \sqrt{\log p/n} \leq \frac{\lambda}{4}$. We, hence, have

$$(a^* - \gamma_0/2) \|\widehat{\nu}\|_2^2 \leq (p_\lambda(\beta_-^*) - p_\lambda(\widehat{\beta}_-)) + \frac{\lambda}{2} \|\widehat{\nu}\|_1. \quad (\text{C.32})$$

(i) Hard sparsity case. Write $\widehat{\nu} = (\widehat{\nu}_1, \widehat{\nu}_-^T)^T$, where $\widehat{\nu}_1 = \widehat{\beta}_1 - \beta_1^*$ corresponding to the intercept term, and $\widehat{\nu}_- = \widehat{\beta}_- - \beta_-^*$. Recall $S_- = \{j : \beta_j^* \neq 0, 2 \leq j \leq p\}$, $S = S_- \cup \{1\}$, and $|S|_0 = s$ is the model sparsity size under the hard sparsity condition. Let $\beta_{S_-}^*$ denote the $(s-1)$ -subvector of β^* consisting of the elements in S_- . Let $\beta_{S_-}^*$, $\widehat{\nu}_{S_-}$, and $\widehat{\nu}_{S^c}$ be defined similarly. Note that under the hard sparsity assumption, $\beta_{S^c}^*$ is a $(p-s)$ -dimensional zero vector. By the subadditivity property of the penalty function $p_\lambda(\cdot)$, we have $p_\lambda(|t_1| + |t_2|) \geq p_\lambda(|t_1|) + p_\lambda(|t_2|)$ for any $t_1, t_2 \in \mathbb{R}$, due to the observation $|t_1| \leq |t_2| + |t_1 + t_2|$. Applying the subadditivity property, we have

$$p_\lambda(\beta_-^*) - p_\lambda(\widehat{\beta}_-) \leq p_\lambda(\beta_{S_-}^*) - [p_\lambda(\beta_{S_-}^* + \widehat{\nu}_{S^c}) - p_\lambda(\widehat{\nu}_{S_-})] = p_\lambda(\widehat{\nu}_{S_-}) - p_\lambda(\widehat{\nu}_{S^c}). \quad (\text{C.33})$$

By the assumption on the penalty function and Lemma 4 of Loh and Wainwright (2015), we have

$$\lambda \|\widehat{v}\|_1 \leq p_\lambda(\widehat{v}) + \frac{\gamma_0}{2} \|\widehat{v}\|_2^2 \leq p_\lambda(\widehat{v}_{S_-}) + p_\lambda(\widehat{v}_{S^c}) + p_\lambda(\widehat{v}_1) + \frac{\gamma_0}{2} \|\widehat{v}\|_2^2. \quad (\text{C.34})$$

Combining (C.32) with (C.33) and (C.34), we have

$$(a^* - 3\gamma_0/4) \|\widehat{v}\|_2^2 \leq \frac{3}{2} p_\lambda(\widehat{v}_{S_-}) - \frac{1}{2} p_\lambda(\widehat{v}_{S^c}) + \frac{1}{2} p_\lambda(\widehat{v}_1) \leq \frac{3}{2} p_\lambda(\widehat{v}_S), \quad (\text{C.35})$$

where $\widehat{v}_S = (\widehat{v}_1, \widehat{v}_{S_-}^T)^T$. As $p_\lambda(t)$ is a λ -Lipschitz continuous function of t (Loh and Wainwright, 2015, Lem. 4), we have $p_\lambda(\widehat{v}_S) \leq \lambda \|\widehat{v}_S\|_1$. Hence,

$$(a^* - 3\gamma_0/4) \|\widehat{v}\|_2^2 \leq \frac{3\lambda}{2} \|\widehat{v}_S\|_1 \leq \frac{3\lambda}{2} \sqrt{s} \|\widehat{v}\|_2.$$

We have proved $\|\widehat{v}\|_2 \leq a_3^* \sqrt{s \log p/n}$, where $a_3^* = \frac{6k_0}{4a^* - 3\gamma_0}$. From the above and the argument for Lemma 1, $\|v_{A^c}\|_1 \leq 3\|v_A\|_1$, where A denotes the index set corresponding to the s -largest (in magnitude) elements of v . Hence, $\|\widehat{v}\|_1 \leq 4a_3^* s \sqrt{\log p/n}$. (ii) Soft sparsity case. Consider $S_a = S_{-a} \cup \{1\}$, where $S_{-a} = \{j : |\beta_j^*| > a, 2 \leq j \leq p\}$ and a is an arbitrary positive constant. Applying the subadditivity property of the penalty function $p_\lambda(\cdot)$, we have

$$\begin{aligned} & p_\lambda(\beta_{-}^*) - p_\lambda(\beta_{-}) \\ & \leq [p_\lambda(\beta_{S_{-a}}^*) + p_\lambda(\beta_{S_a^c}^*)] - [p_\lambda(\beta_{S_{-a}}^* + \widehat{v}_{S_a^c}) - p_\lambda(\beta_{S_a^c}^* + \widehat{v}_{S_{-a}})] \\ & = p_\lambda(\widehat{v}_{S_{-a}}) - p_\lambda(\widehat{v}_{S_a^c}) + p_\lambda(\beta_{S_a^c}^*). \end{aligned} \quad (\text{C.36})$$

Moreover, similarly as (C.34), we have

$$\lambda \|\widehat{v}\|_1 \leq p_\lambda(\widehat{v}) + \frac{\gamma_0}{2} \|\widehat{v}\|_2^2 \leq p_\lambda(\widehat{v}_{S_{-a}}) + p_\lambda(\widehat{v}_{S_a^c}) + p_\lambda(\widehat{v}_1) + \frac{\gamma_0}{2} \|\widehat{v}\|_2^2. \quad (\text{C.37})$$

Combining (C.32) with (C.36) and (C.37) and noting $p_\lambda(\beta_{S_a^c}^*) \leq \lambda \|\beta_{S_a^c}^*\|_1 \leq \lambda R$, we have

$$\begin{aligned} (a^* - 3\gamma_0/4) \|\widehat{v}\|_2^2 & \leq \frac{3}{2} p_\lambda(\widehat{v}_{S_{-a}}) - \frac{1}{2} p_\lambda(\widehat{v}_{S_a^c}) + \frac{1}{2} p_\lambda(\widehat{v}_1) + p_\lambda(\beta_{S_a^c}^*) \\ & \leq \frac{3}{2} p_\lambda(\widehat{v}_{S_{-a}}) - \frac{1}{2} p_\lambda(\widehat{v}_{S_a^c}) + \lambda \|\beta_{S_a^c}^*\|_1. \end{aligned} \quad (\text{C.38})$$

As $p_\lambda(\widehat{v}_{S_a^c}) \geq \lambda \|\widehat{v}_{S_a^c}\|_1 - \frac{\gamma_0}{2} \|\widehat{v}_{S_a^c}\|_2^2$, we have

$$(a^* - \gamma_0) \|\widehat{v}\|_2^2 \leq \frac{3}{2} \lambda \|\widehat{v}_{S_{-a}}\|_1 - \frac{1}{2} \lambda \|\widehat{v}_{S_a^c}\|_1 + \lambda \|\beta_{S_a^c}^*\|_1. \quad (\text{C.39})$$

Similarly, as the derivation for (C.21), $\|\widehat{v}_{S_a}\|_1 \leq \sqrt{a^{-1}R + 1} \|\widehat{v}\|_2$. We have

$$(a^* - \gamma_0) \|\widehat{v}\|_2^2 \leq \frac{3}{2} \sqrt{a^{-1}R + 1} \lambda \|\widehat{v}\|_2 + \lambda R.$$

Hence,

$$\|\widehat{v}\|_2 \leq 2 \max \left\{ \frac{3k_0}{2(a^* - \gamma_0)} \sqrt{a^{-1}R + 1} \sqrt{\log p/n}, \sqrt{k_0/(a^* - \gamma_0)} R^{1/2} (\log p/n)^{1/4} \right\}.$$

As in the proof of Theorem 2, taking $a = \sqrt{\log p/n}$ leads to $\sqrt{a^{-1}R+1}\sqrt{\log p/n} \leq \sqrt{2R(\log p/n)^{1/4}}$. We then obtain

$$\|\widehat{v}\|_2 \leq a_4^* R^{1/2} (\log p/n)^{1/4},$$

where $a_4^* = 2 \max \left\{ \frac{3\sqrt{2}k_0}{2(a^* - \gamma_0)}, \sqrt{\frac{k_0}{a^* - \gamma_0}} \right\}$. On the event of $\widetilde{\Lambda}_n$, the argument for Lemma 1 ensures that $\|\widehat{v}_{S_a^c}\|_1 \leq 3\|\widehat{v}_{S_a}\|_1 + 2\|\beta_{S_a^c}^*\|_1$. Hence,

$$\begin{aligned} \|\widehat{v}\|_1 &\leq 2(2\sqrt{|S_a|}\|\widehat{v}\|_2 + \|\beta_{S_a^c}^*\|_1) \\ &\leq 2(2a_4^*\sqrt{|S_a|}R^{1/2}(\log p/n)^{1/4} + \|\beta_{S_a^c}^*\|_1), \end{aligned}$$

for any $a > 0$, where $|S_a|$ denotes the cardinality of the set S_a . \blacksquare

Proof of Theorem 4. Inequalities (17) and (18) imply that

$$\begin{aligned} |Q_n(\widehat{\beta}) - Q_n(\beta^*)| &\leq \max \{ |S_n(\beta^*)^T (\widehat{\beta} - \beta^*)|, |\bar{S}_n(\widehat{\beta})^T (\widehat{\beta} - \beta^*)| \} \\ &\leq \max \{ \|S_n(\beta^*)\|_\infty, \|\bar{S}_n(\widehat{\beta})\|_\infty \} \|\widehat{\beta} - \beta^*\|_1. \end{aligned}$$

The above inequality holds for any subgradients $S_n(\beta), \bar{S}_n(\beta) \in \partial Q_n(\beta)$. We take $S_n(\beta^*) = n^{-1} \sum_{i=1}^n X_i \xi_i$, where $\xi_i = \mathbf{I}(\epsilon_i < 0) - \tau$. By Lemma 3, $P(\lambda \geq 4\|S_n(\beta^*)\|_\infty) \geq 1 - \delta_n - 2\exp(-\log p)$. We take $\bar{S}_n(\widehat{\beta})$ to be the subgradient satisfying (8), whose existence is guaranteed by the KKT condition for the convex difference program (Tao and An, 1997). Then

$$\bar{S}_n(\widehat{\beta}) + \lambda \text{sgn}(\widehat{\beta}) - H'(\widehat{\beta}_-) = 0, \quad (\text{C.40})$$

where $H'(\cdot)$ is defined in Section 2.3 of the main paper. Note that for QR-LASSO, $H'(\cdot) = 0$ and thus $\|\bar{S}_n(\widehat{\beta})\|_\infty = \lambda$, whereas for QR-NCP, note that $-H(\beta) + \lambda\|\beta\|_1 = p_\lambda(\beta) = \sum_{j=2}^p p_\lambda(|\beta_j|)$. We have $\|H'(\widehat{\beta}) - \lambda \text{sgn}(\widehat{\beta})\|_\infty = \|\partial p_\lambda(\widehat{\beta}_-)\|_\infty$, which is upper bounded by λ (e.g., Loh and Wainwright, 2015, Lem. 4). Hence, $\|\bar{S}_n(\widehat{\beta})\|_\infty \leq \lambda$. Summarizing the above, we have with probability at least $1 - \delta_n - 2\exp(-\log p)$, $|Q_n(\widehat{\beta}) - Q_n(\beta^*)| \leq 4\lambda\|\widehat{\beta} - \beta^*\|_1$. \blacksquare

APPENDIX D. Additional Technical Results

LEMMA D.1. For an arbitrary $k \geq 2$, let $\mathbb{C}(k) = \{\theta \in \mathbb{R}^p : \|\theta\|_2 \leq 1, \|\theta\|_0 \leq k\}$. Then there exist some positive constants α_1 and α_2 such that $\forall t > 0$,

$$\begin{aligned} P \left\{ \sup_{\theta \in \mathbb{C}(k)} |n^{-1}\|\mathbb{X}\theta\|_2^2 - E(n^{-1}\|\mathbb{X}\theta\|_2^2)| \geq t \right\} \\ \leq \alpha_2 \exp \left(-\alpha_1 n \min(t^2/\sigma_x^4, t/\sigma_x^2) + k \log p \right). \end{aligned} \quad (\text{D.1})$$

Proof. This is a minor extension of Lemma 15 of Loh and Wainwright (2012) to allow X to include an intercept term. We provide below an outline of the derivation. First, we show that the exponential inequality in Lemma 14 of Loh and Wainwright (2012) can be extended to allow X to include an intercept term. Specifically, write $\mathbb{X} = (1_n, \widetilde{\mathbb{X}})$, where 1_n denotes an $n \times 1$ column vector of ones and $\widetilde{\mathbb{X}}$ is an $n \times (p-1)$ matrix of covariates where each row is a sub-Gaussian vector and the rows are independent. For an arbitrary $p \times 1$ vector θ , write

$\theta = (\theta_1, \tilde{\theta}')'$. Then $n^{-1} \|\mathbb{X}\theta\|_2^2 = \theta_1^2 + n^{-1} \|\tilde{\mathbb{X}}\tilde{\theta}\|_2^2 + 2\theta_1 n^{-1} 1_n' \tilde{\mathbb{X}}\tilde{\theta}$. For any $\theta \in \mathbb{R}^p$ and any $t > 0$,

$$\begin{aligned} & P\left[\left| \|\mathbb{X}\theta\|_2^2 - E(\|\mathbb{X}\theta\|_2^2) \right| \geq nt \right] \\ & \leq P\left[\left| \|\tilde{\mathbb{X}}\tilde{\theta}\|_2^2 - E(\|\tilde{\mathbb{X}}\tilde{\theta}\|_2^2) \right| \geq nt/2 \right] + P\left[\left| 2\theta_1 n^{-1} 1_n' \tilde{\mathbb{X}}\tilde{\theta} \right| \geq nt/2 \right]. \end{aligned}$$

Consider any fixed θ such that $\|\theta\|_2 \leq 1$. By Lemma 14 of Loh and Wainwright (2012), the first term at the right side of the inequality is upper bounded by $2\exp(-c_1 n \min(t^2/\sigma_x^4, t/\sigma_x^2))$ for some positive constant c_1 . Observing that $n^{-1} 1_n' \tilde{\mathbb{X}}\tilde{\theta}$ is an average of i.i.d. sub-Gaussian random variables, the second term at the right side of the inequality is upper bounded by $\exp(-c_2 nt^2/\sigma_x^2)$ for some positive constant c_2 . Therefore,

$$P\left[\left| \|\mathbb{X}\theta\|_2^2 - E(\|\mathbb{X}\theta\|_2^2) \right| \geq nt \right] \leq c_3 \exp(-c_4 n \min(t^2/\sigma_x^4, t/\sigma_x^2)), \quad (\text{D.2})$$

for some positive constants c_3 and c_4 . Using this exponential inequality and applying the same technique as in the proof of Lemma 15 of Loh and Wainwright (2012) (noting that $s \geq 1$ in that lemma is arbitrary) establishes the desired result. \blacksquare

LEMMA D.2. Assume the conditions of Lemma B.1 are satisfied. We have

$$P(\|\hat{\Sigma}\|_\infty \leq 12\zeta_0^2) \geq 1 - 2\exp(-\log p).$$

Proof. First, by the Cauchy-Schwarz inequality, $|\hat{\Sigma}_{jk}| = |n^{-1} \sum_{i=1}^n x_{ij}x_{ik}| \leq (\hat{\Sigma}_{jj})^{1/2} (\hat{\Sigma}_{kk})^{1/2}$, $\forall 1 \leq j, k \leq p$. Hence, $\|\hat{\Sigma}\|_\infty \leq \max_{1 \leq j \leq p} n^{-1} \sum_{i=1}^n x_{ij}^2$. Since in this lemma x_{ij} is sub-Gaussian with variance proxy bounded by ζ_0^2 , we have $\Sigma_{jj} = E(x_{ij}^2) \leq 4\zeta_0^2$, $j = 1, \dots, p$. A mean-zero random variable x is subexponential with parameter (ζ_*^2, b) , denoted by $\text{SE}(\zeta_*^2, b)$, if $E\{e^{tx}\} \leq \exp(\zeta_*^2 t^2/2)$ for ant $|t| \leq \frac{1}{b}$. The sub-Gaussian property of x_{ij} implies that $x_{ij}^2 - E(x_{ij}^2) \sim \text{SE}(256\zeta_0^4, 16\zeta_0^2)$, $j = 1, \dots, p$.

Applying Bernstein's inequality for subexponential random variables, $\forall t > 0$,

$$P\left(n^{-1} \left| \sum_{i=1}^n (x_{ij}^2 - E(x_{ij}^2)) \right| > t\right) \leq 2\exp\left\{-\frac{n}{2} \min\left(\frac{t^2}{256\zeta_0^4}, \frac{t}{16\zeta_0^2}\right)\right\}. \quad (\text{D.3})$$

Taking $t = 32\zeta_0^2 \sqrt{\log p/n}$ and noting that we assume $\log p \leq n/4$, by the union bound, we have

$$P\left(\max_{1 \leq j \leq p} |\hat{\Sigma}_{jj} - \Sigma_{jj}| > 16\zeta_0^2 \sqrt{\log p/n}\right) \leq 2\exp(-\log p).$$

This implies with probability at least $1 - 2\exp(-\log p)$, we have $P(\|\hat{\Sigma}\|_\infty \leq 12\zeta_0^2)$ \blacksquare

LEMMA D.3. Assume the conditions of Lemma B.1 are satisfied. If $\sqrt{\log p/n} \leq \zeta_1$, where $\zeta_1 = \zeta_*/(32e\zeta_0)$, then

$$P\left(\min_{1 \leq j \leq p} \sum_{i=1}^n |x_{ij}| \geq \zeta_*/2\right) \geq 1 - 2\exp(-\log p).$$

Proof. We will first verify that if a random variable x is sub-Gaussian with variance proxy ζ_0^2 , then $|x|$ is subexponential. The sub-Gaussian property of x implies that $E(|x|^k) \leq (2\zeta_0^2)^{k/2} k\Gamma(k/2)$, for any positive integer $k \geq 1$. We consider below the moment generating function of $|x| - E(|x|)$. For any $t \in \mathbb{R}$,

$$\begin{aligned}
& E\{\exp[t(|x| - E(|x|))]\} \\
&= 1 + \sum_{k=2}^{\infty} \frac{t^k E[(|x| - E(|x|))^k]}{k!} \\
&\leq 1 + \sum_{k=2}^{\infty} \frac{t^k 2^{k-1} E[|x|^k + (E(|x|))^k]}{k!} \\
&\leq 1 + \sum_{k=2}^{\infty} \frac{t^k 2^k E[|x|^k]}{k!} \\
&\leq 1 + \sum_{k=2}^{\infty} \frac{(2t)^k (2\zeta_0^2)^{k/2} k\Gamma(k/2)}{k!} \\
&\leq 1 + \sum_{k=2}^{\infty} (4e\zeta_0 t)^k \\
&\leq 1 + 32e^2 \zeta_0^2 t^2 \quad \text{for } |t| \leq \frac{1}{8e\zeta_0} \\
&\leq \exp(32e^2 \zeta_0^2 t^2),
\end{aligned}$$

where the second last inequality follows by noting $\Gamma(k/2) \leq (k/2)^{k/2}$ and $k^{1/k} \leq e^{1/e}$ for any $k \geq 2$, and by Stirling's approximation $k! \geq (k/e)^k$. Hence, $|x| - E(|x|) \sim \text{SE}(64e^2 \zeta_0^2, 8e\zeta_0)$. Applying the Bernstein's inequality for subexponential random variables, we have

$$P\left(n^{-1} \left| \sum_{i=1}^n (|x_{ij}| - E(|x_{ij}|)) \right| > t\right) \leq 2 \exp \left\{ -\frac{n}{2} \min \left(\frac{t^2}{64e^2 \zeta_0^2}, \frac{t}{8e\zeta_0} \right) \right\}, \quad (\text{D.4})$$

$\forall t > 0$. Taking $t = 16e\zeta_0 \sqrt{\log p/n}$ in (D.4) and noting that we assume $\log p \leq n/4$, by the union bound, we have

$$P\left(\max_{1 \leq j \leq p} n^{-1} \left| \sum_{i=1}^n (|x_{ij}| - E(|x_{ij}|)) \right| > 16e\zeta_0 \sqrt{\log p/n}\right) \leq 2 \exp(-\log p).$$

Hence, with probability at least $1 - 2 \exp(-\log p)$,

$$\sum_{i=1}^n |x_{ij}| \geq E(|x_{ij}|) - 16e\zeta_0 \sqrt{\log p/n}, \text{ for all } j = 1, \dots, p.$$

As $\min_{1 \leq j \leq p} E(|x_{ij}|) \geq \zeta_* > 0$ and $32e\zeta_0 \sqrt{\log p/n} \leq \zeta_*$, the conclusion of the lemma follows. \blacksquare

Proof of Lemma 1. This result follows from the generalized KKT condition of the convex difference program (Tao and An, 1997). We provide a self-contained derivation below.

(i) The result for the global minimum of QR-LASSO follows directly from the definition of subgradient. To see this, let $\hat{\beta} = \arg \min_{\beta} \{Q_n(\beta) + \lambda|\beta_-|\}$ denote the global minimum of QR-LASSO. Then

$$\{Q_n(\beta) + \lambda|\beta_-|\} - \{Q_n(\hat{\beta}) + \lambda|\hat{\beta}_-|\} \geq 0_p^T(\beta - \hat{\beta}) = 0,$$

where 0_p denotes a p -dimensional zero vector. Hence, by the definition of subgradient, we have $0_p \in \partial\{Q_n(\hat{\beta}) + \lambda|\hat{\beta}_-|\}$.

(ii) Now, consider the case of QR-NCP. Let $\hat{\beta}$ denote a local minimum of $L_n(\beta) = \tilde{L}_n(\beta) - H(\beta)$, where $\tilde{L}_n(\beta) = Q_n(\beta) + \lambda \sum_{j=2}^p |\beta_j|$ and $H(\beta) = \sum_{j=2}^p h_{\lambda}(\beta_j)$. Then there exists a neighborhood U of $\hat{\beta}$ such that $L_n(\beta) \geq L_n(\hat{\beta})$, for all $\beta \in U$. Hence, $\forall \beta \in U$,

$$\tilde{L}_n(\beta) - \tilde{L}_n(\hat{\beta}) \geq H(\beta) - H(\hat{\beta}) \geq (H'(\hat{\beta}))^T(\beta - \hat{\beta}), \quad (\text{D.5})$$

where $H'(\hat{\beta}) = (0, h'_{\lambda}(\hat{\beta}_2), \dots, h'_{\lambda}(\hat{\beta}_p))^T$, and the second inequality follows because $H(\beta)$ is convex and differentiable. The convexity of L_n and (D.5) implies that $H'(\hat{\beta}) \in \partial\tilde{L}_n(\hat{\beta})$. This finishes the proof. ■

Proof of Lemma 2 (QR-LASSO). The proof of the result under hard sparsity was given in Belloni and Chernozhukov (2011). We include it here for completeness and to facilitate the proof for the result under soft sparsity. Consider the event $\Lambda_n = \{\lambda \geq 2\|\tilde{S}_n\|_{\infty}\}$, where λ satisfies the conditions of Lemma 2. Lemma 3 ensures that $P(\Lambda_n) \geq 1 - \delta_n - 2\exp(-\log p)$. Recall $\hat{v} = \hat{\beta} - \beta^*$. By the definition of $\hat{\beta}$, $Q_n(\hat{\beta}) + \lambda\|\hat{\beta}_- \|_1 \leq Q_n(\beta^*) + \lambda\|\beta^*_- \|_1$. This implies

$$Q_n(\hat{\beta}) - Q_n(\beta^*) \leq \lambda(\|\beta^*_- \|_1 - \|\hat{\beta}_- \|_1) \leq \lambda(\|\hat{v}_S \|_1 + \|\hat{v}_{S^c} \|_1). \quad (\text{D.6})$$

On the other hand, the convexity of $Q_n(\cdot)$ guarantees that on Λ_n ,

$$Q_n(\hat{\beta}) - Q_n(\beta^*) \geq \tilde{S}_n \hat{v} \geq -\|\hat{v}\|_1 \|\tilde{S}_n\|_{\infty} \geq -\frac{\lambda}{2}(\|\hat{v}_S \|_1 + \|\hat{v}_{S^c} \|_1). \quad (\text{D.7})$$

Putting (D.6) and (D.7) together, we have $\hat{v} \in \Gamma_H$.

Under the soft sparsity assumption, by a similar argument as above for (D.6) and (D.7), we obtain that on Λ_n ,

$$2(\|\beta^*_- + \hat{v}_- \|_1 - \|\beta^*_- \|_1) \leq \|\hat{v}_{S_a} \|_1 + \|\hat{v}_{S_a^c} \|_1, \quad (\text{D.8})$$

where $S_{-a} = \{j : |\beta_j^*| > a, 2 \leq j \leq p\}$ and $S_a = S_{-a} \cup \{1\}$. On the other hand,

$$\begin{aligned} & \|\beta^*_- + \hat{v}_- \|_1 - \|\beta^*_- \|_1 \\ & \geq (\|\beta_{S_{-a}}^* + \hat{v}_{S_a^c} \|_1 - \|\beta_{S_a^c}^* + \hat{v}_{S_{-a}} \|_1) - (\|\beta_{S_{-a}}^* \|_1 + \|\beta_{S_a^c}^* \|_1) \\ & = \|\hat{v}_{S_a^c} \|_1 - \|\hat{v}_{S_{-a}} \|_1 - 2\|\beta_{S_a^c}^* \|_1 \\ & \geq \|\hat{v}_{S_a^c} \|_1 - \|\hat{v}_{S_a} \|_1 - 2\|\beta_{S_a^c}^* \|_1. \end{aligned} \quad (\text{D.9})$$

Combining (D.8) with (D.9), we have

$$2(\|\hat{v}_{S_a^c} \|_1 - \|\hat{v}_{S_a} \|_1 - 2\|\beta_{S_a^c}^* \|_1) \leq \|\hat{v}_{S_a} \|_1 + \|\hat{v}_{S_a^c} \|_1.$$

Hence, $\hat{v} \in \Gamma_W$ under the soft sparsity assumption. ■

Proof of Lemma 3. Let $A_{n1} = \{\max_{1 \leq j \leq p} \hat{\sigma}_j^2 \leq m_x\}$, where $\hat{\sigma}_j^2 = n^{-1} \sum_{i=1}^n x_{ij}^2$. Then $P(A_{n1}) \geq 1 - \delta_n$, where $\delta_n \rightarrow 0$ as $n \rightarrow \infty$ by condition (C3). On the event A_{n1} , by the union bound, we have

$$\begin{aligned} P(\|\tilde{S}_n\|_\infty > 2\sqrt{m_x \log p/n}) &= P\left(\max_{1 \leq j \leq p} \left|n^{-1} \sum_{i=1}^n x_{ij} \xi_i\right| > 2\sqrt{m_x \log p/n}\right) \\ &\leq \sum_{j=1}^p P\left(\left|n^{-1} \sum_{i=1}^n x_{ij} \xi_i\right| > 2\sqrt{m_x \log p/n}\right). \end{aligned}$$

As $-\tau \leq \xi_i \leq 1 - \tau$, ξ is a sub-Gaussian random variable with parameter bounded by one. Hence,

$$\begin{aligned} P\left(\left|n^{-1} \sum_{i=1}^n x_{ij} \xi_i\right| > 2\sqrt{m_x \log p/n}\right) \\ \leq 2E_X \left\{ \exp\left(-\frac{4m_x \log p}{2n^{-1} \sum_{i=1}^n x_{ij}^2}\right) \right\} \leq 2\exp(-2\log p). \end{aligned}$$

We have

$$P(\|\tilde{S}_n\|_\infty > 2\sqrt{m_x \log p/n}) \leq 2\exp(\log p - 2\log p) = 2\exp(-\log p).$$

This proves the lemma. ■

Proof of Corollary 1 (QR-NCP). The geometric structures of the local solutions for QR-NCP are implied by the derivation in the proof of Theorem 3.

For the hard sparsity case, it follows from the first inequality of (C.35) that

$$\begin{aligned} (a^* - 3\gamma_0/4) \|\hat{v}\|_2^2 &\leq \frac{3}{2} p_\lambda(\hat{v}_{S_-}) - \frac{1}{2} p_\lambda(\hat{v}_{S^c}) + \frac{1}{2} p_\lambda(\hat{v}_1) \\ &\leq \frac{3}{2} p_\lambda(\hat{v}_S) - \frac{1}{2} p_\lambda(\hat{v}_{S^c}) \\ &\leq \frac{3}{2} p_\lambda(\hat{v}_A) - \frac{1}{2} p_\lambda(\hat{v}_{A^c}), \end{aligned}$$

where A is the index set of the largest s elements of \hat{v} in magnitude. This implies $3p_\lambda(\hat{v}_A) - p_\lambda(\hat{v}_{A^c}) \geq 0$. Lemma 5 of Loh and Wainwright (2015) implies that $0 \leq 3p_\lambda(\hat{v}_A) - p_\lambda(\hat{v}_{A^c}) \leq \lambda(3\|\hat{v}_A\|_1 - \|\hat{v}_{A^c}\|_1)$ or $\|\hat{v}_{A^c}\|_1 \leq 3\|\hat{v}_A\|_1$.

For the soft sparsity case, it follows from (C.39) that

$$\frac{1}{2} \lambda \|\hat{v}_{S_a^c}\|_1 \leq \frac{3}{2} \lambda \|\hat{v}_{S_a}\|_1 + \lambda \|\beta_{S_a^c}^*\|_1,$$

$$\text{or } \|\hat{v}_{S_a^c}\|_1 \leq 3\|\hat{v}_{S_a}\|_1 + 2\|\beta_{S_a^c}^*\|_1. \quad \blacksquare$$

Proof of Corollary 2. The result in (i) for QR-NCP follows immediately from Theorem 4. To establish the result in (i) for QR-LASSO, let $\hat{v} = \hat{\beta} - \beta^* = (\hat{v}_1, \hat{v}_-)'$, where $\hat{\beta}$ denotes the QR-LASSO estimator. By the definition of $\hat{\beta}$, we have

$$|R_n(\hat{\beta})| \leq \lambda(\|\beta_-^*\|_1 - \|\hat{v}_-\|_1). \quad (\text{D.10})$$

Note that as $\|\beta_-^* + \widehat{v}_- \|_1 \geq \|\beta_-^* \|_1 - \|\widehat{v}_- \|_1$, the right-hand side of (D.10) immediately implies that

$$|R_n(\widehat{\beta})| \leq \lambda \left(2\|\beta_-^* \|_1 - \|\widehat{v}_- \|_1 \right) \leq 2\lambda \|\beta_-^* \|_1 \leq 2\lambda \|\beta^* \|_1.$$

The results in (ii) and (iii) follow immediately by combining Theorem 4 with the L_1 -estimation error bound derived in Theorems 2 and 3. \blacksquare

REFERENCES

- Abadie, A., J. Angrist, & G. Imbens (2002) Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings. *Econometrica* 70(1), 91–117.
- Angrist, J., V. Chernozhukov, & I. Fernández-Val (2006) Quantile regression under misspecification, with an application to the US wage structure. *Econometrica* 74(2), 539–563.
- Arellano, M. & S. Bonhomme (2017) Quantile selection models with an application to understanding changes in wage inequality. *Econometrica* 85(1), 1–28.
- Belloni, A. & V. Chernozhukov (2011) L1-penalized quantile regression in high-dimensional sparse models. *Annals of Statistics* 39, 82–130.
- Belloni, A., V. Chernozhukov, & K. Kato (2014) Uniform post-selection inference for least absolute deviation regression and other z-estimation problems. *Biometrika* 102(1), 77–94.
- Belloni, A., V. Chernozhukov, & K. Kato (2019) Valid post-selection inference in high-dimensional approximately sparse quantile regression models. *Journal of the American Statistical Association* 114(526), 749–758.
- Bickel, P.J., Y. Ritov, & A.B. Tsybakov (2009) Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics* 37(4), 1705–1732.
- Bradic, J., J. Fan, & W. Wang (2011) Penalized composite quasi-likelihood for ultrahigh dimensional variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(3), 325–349.
- Bradic, J. & M. Kolar (2017). Uniform inference for high-dimensional quantile regression: Linear functionals and regression rank scores. *Preprint*, [arXiv:1702.06209](https://arxiv.org/abs/1702.06209).
- Buchinsky, M. (1994) Changes in the US wage structure 1963–1987: Application of quantile regression. *Econometrica* 62, 405–458.
- Buchinsky, M. (1998) The dynamics of changes in the female wage distribution in the USA: A quantile regression approach. *Journal of Applied Econometrics* 13(1), 1–30.
- Bunea, F., A. Tsybakov, & M. Wegkamp (2007) Sparsity oracle inequalities for the Lasso. *Electronic Journal of Statistics* 1, 169–194.
- Chamberlain, G. (1994) Quantile regression, censoring, and the structure of wages. In C.A. Sims (ed.), *Advances in Econometrics: Sixth World Congress*, vol. 2. Cambridge University Press, pp. 171–209.
- Chen, X., D. Li, Q. Li, & Z. Li (2019a) Nonparametric estimation of conditional quantile functions in the presence of irrelevant covariates. *Journal of Econometrics* 212(2), 433–450.
- Chen, X., W. Liu, & Y. Zhang (2019b) Quantile regression under memory constraint. *Annals of Statistics* 47(6), 3244–3273.
- Chernozhukov, V. & I. Fernández-Val (2011) Inference for extremal conditional quantile models, with an application to market and birthweight risks. *The Review of Economic Studies* 78(2), 559–589.
- Chernozhukov, V., I. Fernández-Val, J. Hahn, & W. Newey (2013) Average and quantile effects in nonseparable panel models. *Econometrica* 81(2), 535–580.
- Donoho, D.L. & I.M. Johnstone (1994). Minimax risk over l_p -balls for l_q -error. *Probability Theory and Related Fields* 99(2), 277–303.
- Elsener, A. & S. van de Geer (2018) Sharp oracle inequalities for stationary points of nonconvex penalized m -estimators. *IEEE Transactions on Information Theory* 65(3), 1452–1472.

- Fan, J., Y. Fan, & E. Barut (2014) Adaptive robust variable selection. *Annals of Statistics* 42(1), 324–351.
- Fan, J., Q. Li, & Y. Wang (2017) Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79(1), 247–265.
- Fan, J. & R. Li (2001) Variable selection via nonconcave penalized likelihood and its oracle property. *Journal of the American Statistical Association* 96, 1348–1360.
- Fan, Z. & H. Lian (2018) Quantile regression for additive coefficient models in high dimensions. *Journal of Multivariate Analysis* 164, 54–64.
- Firpo, S., N.M. Fortin, & T. Lemieux (2009) Unconditional quantile regressions. *Econometrica* 77(3), 953–973.
- Fitzenberger, B., R. Koenker, & J.A. Machado (2013) *Economic Applications of Quantile Regression*. Springer Science & Business Media.
- Galvao, A.F., C. Lamarche, & L.R. Lima (2013) Estimation of censored quantile regression for panel data with fixed effects. *Journal of the American Statistical Association* 108(503), 1075–1089.
- Graham, B.S., J. Hahn, A. Poirier, & J.L. Powell (2018) A quantile correlated random coefficients panel data model. *Journal of Econometrics* 206(2), 305–335.
- Greenshtein, E., Y. Ritov (2004) Persistence in high-dimensional linear predictor selection and the virtue of overparametrization. *Bernoulli* 10(6), 971–988.
- Harding, M. & C. Lamarche (2018) A panel quantile approach to attrition bias in big data: Evidence from a randomized experiment. *Journal of Econometrics* 211, 61–82.
- Honda, T., C.-K. Ing, & W.-Y. Wu (2019) Adaptively weighted group lasso for semiparametric quantile regression models. *Bernoulli* 25(4B), 3311–3338.
- Horowitz, J.L. & S. Lee (2005) Nonparametric estimation of an additive quantile regression model. *Journal of the American Statistical Association* 100(472), 1238–1249.
- Horowitz, J.L. & V.G. Spokoiny (2002) An adaptive, rate-optimal test of linearity for median regression models. *Journal of the American Statistical Association* 97(459), 822–835.
- Kai, B., R. Li, & H. Zou (2011) New efficient estimation and variable selection methods for semiparametric varying-coefficient partially linear models. *Annals of Statistics* 39, 305–332.
- Kato, K. (2011) Group Lasso for high dimensional sparse quantile regression models. *Preprint*, [arXiv:1103.1458](https://arxiv.org/abs/1103.1458).
- Koenker, R. (2017) Quantile regression: 40 years on. *Annual Review of Economics* 9, 155–176.
- Koenker, R. & G. Bassett (1978) Regression quantiles. *Econometrica* 46, 33–50.
- Koenker, R., V. Chernozhukov, X. He, & L. Peng (eds.) (2017) *Handbook of Quantile Regression*. Chapman and Hall/CRC.
- Koenker, R. & Z. Xiao (2006) Quantile autoregression. *Journal of the American Statistical Association* 101(475), 980–990.
- Ledoux, M. & M. Talagrand (2013) *Probability in Banach Spaces: Isoperimetry and Processes*. Springer Science & Business Media.
- Lee, E.R., H. Noh, & B.U. Park (2014) Model selection via Bayesian information criterion for quantile regression models. *Journal of the American Statistical Association* 109(505), 216–229.
- Lee, S., Y. Liao, M.H. Seo, & Y. Shin (2018) Oracle estimation of a change point in high dimensional quantile regression. *Journal of the American Statistical Association* 43, 1184–1194.
- Li, Y.J. & J. Zhu (2008) L1-norm quantile regression. *Journal of Computational and Graphical Statistics* 17, 163–185.
- Linton, O.B. & Y.-J. Whang (2004). A quantilogram approach to evaluating directional predictability. Available at SSRN 485342.
- Loh, P.-L. (2017). Statistical consistency and asymptotic normality for high-dimensional robust m -estimators. *Annals of Statistics* 45(2), 866–896.
- Loh, P.-L. & M.J. Wainwright (2012) High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *Annals of Statistics* 40(3), 1637–1664.

- Loh, P.-L. and M.J. Wainwright (2015). Regularized m -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research* 16, 559–616.
- Lv, S., H. Lin, H. Lian, & J. Huang (2018) Oracle inequalities for sparse additive quantile regression in reproducing kernel Hilbert space. *Annals of Statistics* 46(2), 781–813.
- Mei, S., Y. Bai, & A. Montanari (2018) The landscape of empirical risk for nonconvex losses. *Annals of Statistics* 46(6A), 2747–2774.
- Negahban, S.N., P. Ravikumar, M.J. Wainwright, & B. Yu (2012) A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. *Statistical Science* 27(4), 538–557.
- Nolan, J. (2003) *Stable Distributions: Models for Heavy-Tailed Data*. Birkhauser.
- Park, S., X. He, & S. Zhou (2017) Dantzig-type penalization for multiple quantile regression with high dimensional covariates. *Statistica Sinica* 27, 1619–1638.
- Raskutti, G., M.J. Wainwright, & B. Yu (2011) Minimax rates of estimation for high-dimensional linear regression over l_q -balls. *IEEE Transactions on Information Theory* 57(10), 6976–6994.
- Ruppert, D. & R.J. Carroll (1980) Trimmed least squares estimation in the linear model. *Journal of the American Statistical Association* 75(372), 828–838.
- Sherwood, B. & L. Wang (2016) Partially linear additive quantile regression in ultra-high dimension. *Annals of Statistics* 44(1), 288–317.
- Shows, J.H., W. Lu, & H.H. Zhang (2010) Sparse estimation and inference for censored median regression. *Journal of Statistical Planning and Inference* 140, 1903–1917.
- Su, L. & T. Hoshino (2016) Sieve instrumental variable quantile regression estimation of functional coefficient models. *Journal of Econometrics* 191(1), 231–254.
- Tang, Y., X. Song, H.J. Wang, & Z. Zhu (2013) Variable selection in high-dimensional quantile varying coefficient models. *Journal of Multivariate Analysis* 122, 115–132.
- Tao, P.D. & L. An (1997) Convex analysis approach to D.C. programming: Theory, algorithms and applications. *Acta Mathematica Vietnamica* 22(1), 289–355.
- Tibshirani, R. (1996) Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B* 58, 267–288.
- van de Geer, S.A. (2000) *Empirical Processes in M-Estimation*. Cambridge University Press.
- van de Geer, S.A. (2016) *Estimation and Testing under Sparsity*. Springer.
- van der Vaart, A. & J. Wellner (1996) *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Science & Business Media.
- Wagener, J., S. Volgushev, & H. Dette (2012) The quantile process under random censoring. *Mathematical Methods of Statistics* 21, 127–141.
- Wang, H., G. Li, & G. Jiang (2007) Robust regression shrinkage and consistent variable selection through the LAD-lasso. *Journal of Business & Economic Statistics* 25, 347–355.
- Wang, H., J. Zhou, & Y. Li (2013a) Variable selection for censored quantile regression. *Statistica Sinica* 23, 145–167.
- Wang, L. (2013) The L1 penalized LAD estimator for high dimensional linear regression. *Journal of Multivariate Analysis* 120, 135–151.
- Wang, L. (2019). L_1 -regularized quantile regression with many regressors under lean assumptions. *University of Minnesota Digital Conservancy*. Available at <https://hdl.handle.net/11299/202063>.
- Wang, L., Y. Kim, & R. Li (2013b) Calibrating non-convex penalized regression in ultra-high dimension. *Annals of Statistics* 41(5), 2505–2536.
- Wang, L., B. Peng, J. Bradic, R. Li, & Y. Wu (2020) A tuning-free robust and efficient approach to high-dimensional regression. *Journal of the American Statistical Association* 115(532), 1700–1714.
- Wang, L., Y. Wu, & R. Li (2012) Quantile regression for analyzing heterogeneity in ultra-high dimension. *Journal of the American Statistical Association* 107(497), 214–222.
- Wu, Y.C. & Y.F. Liu (2009) Variable selection in quantile regression. *Statistica Sinica* 19, 801–817.
- Zhang, C.H. (2010) Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics* 38, 894–942.

- Zhao, T., M. Kolar, & H. Liu (2014) A general framework for robust testing and confidence regions in high-dimensional quantile regression. *Preprint*, [arXiv:1412.8724](#).
- Zheng, Q., L. Peng, & X. He (2015) Globally adaptive quantile regression with ultra-high dimensional data. *Annals of Statistics* 43(5), 2225–2258.
- Zhong, W., L. Zhu, R. Li, & H. Cui (2016) Regularized quantile regression and robust feature screening for single index models. *Statistica Sinica* 26(1), 69–95.
- Zou, H. & M. Yuan (2008) Composite quantile regression and the oracle model selection theory. *Annals of Statistics* 36, 1108–1126.